

nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

MIND GAMES

What can modern neuroscience tell us about
human brain function? **PAGE 371**

BEHAVIOUR

GUT FEELINGS

Are intestinal bacteria
shaping our minds?

PAGE 312

POLICYMAKING

HOW TO USE EXPERTS

There's a science to
taking science advice

PAGE 317

NEURONAL DEVELOPMENT

IT'S DIFFERENT FOR MALES

Sex-specific neurons linked
to learning in the roundworm

PAGE 385

NATUREASIA.COM

18 October 2015

Vol 526, No 7571

THIS WEEK

EDITORIALS

NEUROSCIENCE The extra neurons on a male worm's mind **p.294**

WORLD VIEW Look more closely at the benefits of alternative medicine **p.295**



FOSSILS Handy find shows that early humans hung in trees **p.297**

A shift in climate

The Intergovernmental Panel on Climate Change has done much to alert politicians to the effects of global warming. But to push climate change up the agenda, it will need to do the same for the public.

Hoesung Lee has laid out a vision for his tenure as the fourth chair of the Intergovernmental Panel on Climate Change (IPCC). The South Korean economist says that he wants to increase coordination between the working groups, work with more scientists from developing countries, boost interaction with the business and financial industries and expand the panel's influence by making its findings easier to digest. Above all, Lee says that he wants to be remembered as the man who shifted the panel's focus towards solutions. Those are all worthy goals, but the organization that assesses — and in some senses oversees — the world of climate science is well placed to do a lot more.

The basic science underlying climate change has been firmly established, and we now know much more about not only the threats posed by a rapidly warming world, but also the options humanity has for changing course. Even though the commitments that have been proffered going into the United Nations climate summit in Paris at the end of next month are generally tepid, the governments of the world by and large recognize the problem and know that they need to act. What else can the IPCC do? The answer is plenty, and Lee's emphasis on solutions is one piece of the puzzle.

Most of the world's major emitters — rich, poor and in between — have offered commitments and policies to move their countries in the right direction. Without concerted action in the decades to come, however, these commitments will fall well short of what the best science suggests is needed to achieve the formal goal of limiting the rise in global average temperatures to 2°C. Political leaders are leaving the hard work for later, and could well be committing future generations to more warming than anybody wants to experience.

This is partly because of a lack of political leadership and of active opposition by entrenched industrial interests, as environmentalists argue. But it is also evidence of the vastness of the challenge. It is not easy to transform the global economy and industrial base for a large and growing population, much of which is still mired in poverty but wants access to the modern conveniences that so many on the planet take for granted. The world needs a full suite of technologies that are not only cost-effective but also socially and politically viable. Here, the IPCC has a particular part to play.

The weak commitments going into Paris are also evidence of a disconnect between scientists, who think that the evidence speaks for itself, and citizens and policymakers, who have a lot of other things on their minds. The IPCC is looking at bringing science writers and graphics experts on board in an effort to improve its reports. A linguistics study published earlier this month showed that the IPCC summaries for policymakers score low in terms of readability, and recommends that key panel members receive science-communication training (R. Barkemeyer *et al. Nature Clim. Change* <http://doi.org/79f>; 2015). All this makes sense, but communicating the science more clearly is just the first, and a relatively minor, step.

The IPCC's reports are aimed mainly at — and written in coordination with — governments, yet politicians at the very highest levels are already talking about climate change. The unfortunate truth is that taking steps to combat climate change is way down the political agenda, and that makes more aggressive action difficult. The real challenge is to raise public awareness about the risks of inaction — as well as the benefits of action — and to identify policies that can pass the political litmus test.

Here, as Lee himself has said, the IPCC has an important role. The panel must generate and incorporate knowledge about how information filters through society and about the

"The real challenge is to raise public awareness about the risks of inaction."

kinds of policies that are most likely to work. This is the domain of sociologists, psychologists, anthropologists and political scientists, and they must be an integral part of the IPCC's sixth assessment.

The IPCC has had its controversies, including a glitch in its 2007 projection for Himalayan glacier melt and this year's resignation of former chairman Rajendra Pachauri, who faces — and denies — accusations of sexual harassment. But the challenges that face the panel today are in many ways a result of its success. Much ground has been covered; the challenge now, for both researchers and the IPCC, is to adapt and to identify research that will help policymakers to bridge the gap between what they say they want to do and what they are actually doing. ■

After Asilomar

Scientist-led conferences are no longer the best way to resolve debates on controversial research.

In 1975, some 140 scientists met at the Asilomar resort on California's rocky Monterey Peninsula to discuss the nascent science of mixing DNA from different organisms.

Until that point, researchers had deliberately not performed the final steps of such experiments, owing to concerns about safety and ethics. Over three days of discussions, the conference attendees agreed to voluntary restrictions on recombinant-DNA research, and drafted a document that listed the potential risks of such experiments and how to carry out the work safely.

The meeting is seen as the first time that science had regulated itself — effectively avoiding government intervention — and assuaged public fears by addressing biosafety concerns head-on.

Today, no scientific controversy is complete without calls for an

'Asilomar-like' conference. Until such a conversation has taken place, proponents say, researchers should not proceed with risky propositions.

Debates on artificial intelligence, autonomous weapons, geoengineering and the use of gene-editing technology have all referred to Asilomar as a useful model. (Geoengineers went so far as to meet in Asilomar.) This month, a group of scholars, programmers, artists, entrepreneurs and video-game developers published a Biosphere Code to protect people and the planet from the negative impact of computer algorithms — produced after Asilomar-like discussions in Stockholm.

But is Asilomar's reputation deserved? The invitation-only conference included a handful of journalists and policymakers, but did not cast a wide net outside the scientific community. And in hindsight, many of its safety precautions may have been overkill. As bioethicist Jonathan Moreno puts it: "Asilomar has become a bio-Woodstock in people's memories, a golden age. People forget how muddy Woodstock was."

Modern science is muddier still: in a 2008 essay in *Nature*, even Asilomar organizer Paul Berg admitted that such a conference would be difficult to convene today (P. Berg *Nature* 455, 290–291; 2008). In 1975, he and his colleagues had yielded to concerns from within their tight-knit community. They could afford to pause their research, having reasonable certainty that the technology would not advance in the meantime.

But like everything else in the twenty-first century, science has become a global affair. An enormous number of researchers have almost unfettered access to information and increasingly easy-to-use tools. As a result, 'synthetic' organisms, enhanced influenza viruses and genetically modified human embryos already exist, whether the world is ready for them or not. Even if they are destroyed, the instructions to make them will inevitably make their way onto the Internet — a technology as pervasive and uncontrollable as any biological entity.

Modern science is also less insular than that of the past, and any single Asilomar conference would probably be lost in the noise. New players have appeared over the past four decades, including a powerful biotechnology industry driven at least in part by profit; the most polarized US government in history, which can turn any new technology into a political weapon; and a mass of religious and activist groups that

have flexed their muscles to stop research on embryonic stem cells and genetically modified organisms in their tracks.

Each represents and interacts separately with the general public. And scientists who wish to self-regulate ignore public outcry at their peril: crowd-pleasing politicians passing knee-jerk regulations will hinder scientific progress more than any voluntary moratorium ever could, and their poor understanding might cause collateral damage to related fields.

"When controversy comes calling, scientists should reach outwards."

When controversy comes calling, rather than asking for an Asilomar conference — which, after all, was closed to the public — scientists should reach outwards. Discussions should extend beyond researchers and ethicists to include, or at least broadcast to, the broader public. Proactive engagement with the mass media is key: the most transparent of webcasts is meaningless if only a rarified group of already-interested individuals knows that the meeting is happening.

The most important thing is to communicate the risks and benefits of controversial research in a responsible and transparent way. For embryo editing, for instance, discussions should avoid unhelpful references to the genetically modified humans in the 1997 film *Gattaca*, or veiled aspersions on the ethical standards of non-Western researchers.

A modern Asilomar might also take advantage of the wide range of expertise and techniques available today. Some advisory bodies, such as the US National Academies' committee on research that involves enhancing influenza viruses, have had the forethought to include economists and futurists expert in drawing up realistic risk–benefit analyses and scenarios. Other strategies could include a war-game approach similar to the 2001 Dark Winter exercise, in which US media, government officials, health experts and military groups simulated a bioterror attack to anticipate the problems that would arise.

In 1975, the week after the Asilomar conference, actor Telly Savalas — star of the detective programme *Kojak* — topped the UK music charts with a spoken-word version of *If* by rock group Bread. The world has moved on since then; science must as well. ■

The worm returns

The wiring diagram of the male nematode's nervous system is only a beginning.

When scientists seemed to have completed the map of the nervous system of a tiny male worm in 2012 (T. A. Jarrell *et al. Science* 337, 437–444; 2012), some researchers were already questioning whether the whole effort, originating some 40 years before, was truly worth it. The construction of the wiring diagram for the nervous system of the male of the nematode species *Caenorhabditis elegans* built on the wiring diagram for the hermaphrodite, established more than 25 years earlier, and required painstaking tracing of how the male worm's extra neurons connected to each other.

Stephen Hawking was talking metaphorically when he famously wrote that to unravel the laws of nature would be to know the mind of God. The *C. elegans* project was quite literal: as Sydney Brenner, the originator of the project, jokingly entitled the manuscript of a landmark 1984 paper, the scientists really did want to know 'the mind of the worm'. And in doing so, they argued, they could learn more about how brains create behaviour in higher organisms all the way up to humans.

Can we really know the mind of a worm? Three years on we have an answer of sorts: possibly. In fact, it turns out that we did not even find all the neurons that comprise the male worm's mind. For on page 385, researchers including two of the 2012 team publish an appendix to

the wiring diagram of the male *C. elegans*. As well as the 383 neurons already identified, they describe the discovery of neurons number 384 and 385, which they found in the worm's head.

There is much to admire about the new work, not least that the researchers chose to call the new cells — mystery cells found in the male worms — MCMs. No metaphor there either. It stands for mystery cells of the male.

What's on a male nematode's MCMs? Not so much of a mystery as it turns out: sex. The new neurons have an old role, and help the worms to learn to prioritize the search for a mate over the need for food. When these neurons are put out of action, the male worms never discover the facts of life. The findings offer much more than a completion of the neural map of the male worm. In most organisms, sex-specific differences in behaviour extend to cognitive-like processes such as learning, which can aid reproductive success, but the underlying neural mechanisms are mostly unclear.

The discovery of the MCMs, and the subsequent experiments with them, link developmental and anatomical sex differences in high-order processing areas of the brain to sex-specific behaviour during learning. In doing so, they help to shed light on the neural basis of sex differences in behaviour. And they show that these neurons arise upon sexual maturation from specialized cells called glia — unlike other neurons in *C. elegans* and other invertebrates, which arise from epithelial or undifferentiated blast cells.

Was the effort to trace out the connections between the male nematode's neurons worth it? Like all good maps, the wiring diagram of the *C. elegans* is best viewed as a starting point. The final destination is sure to surprise us. ■

➔ **NATURE.COM**
To comment online,
click on Editorials at:
go.nature.com/xhunqv

GARRY SIMPSON



Consider all the evidence on alternative therapies

Investigate and incorporate the mechanisms of complementary medicine instead of rejecting it outright, says Jo Marchant.

“Insane”, “a joke”, and “exactly the sort of thing the NHS should not be doing!” are a few of the Twitter responses to last week’s news that Britain’s Princess Alexandra Hospital NHS Trust wants to hire a reiki therapist for a hospital in Epping. On a salary of up to £22,236 (US\$34,000) a year, the appointed person “will provide Reiki/Spiritual healing to patients to enable them to cope with the emotional, physical and spiritual issues of dealing with their cancer journey”.

Critics of the advert — and there are many — advocate instead what they call “evidence-based” approaches to health care. These critics should look again at the evidence — because it shows that to dismiss the benefits of alternative therapies is simplistic and misguided.

Let’s be clear, I don’t buy into the pseudoscientific claims of reiki and spiritual healers. There is no evidence that they can tap into and manipulate human ‘energy fields’ to clear blockages and heal the body. Like many alternative therapies, these practices perform no better than placebos in clinical trials.

But that does not mean that such treatments have no distinct therapeutic value. To dismiss people’s complex psychological and physiological reactions to serious illness — and how it is treated — as mere placebo effects is not helpful.

Neuroscience studies show that placebo effects can trigger significant physiological responses that are often identical to those created by drugs, ranging from the release of dopamine in the brains of people with Parkinson’s disease to a rush of endorphins for those in pain.

The standard ‘evidence-based’ argument is that this is irrelevant. Even if alternative therapies induce a biological response, sceptics argue, patients are still better off receiving trial-proven conventional treatments, because then they benefit from both a placebo effect and the active effect of the drug.

This logic misunderstands the nature of placebo effects. Not all placebos are the same, and alternative therapies can sometimes trigger larger responses than conventional ones do. For example, in one trial, fake acupuncture relieved pain more effectively than a fake pill (T. J. Kaptchuk *et al. Br. Med. J.* **332**, 391–397; 2006); in another, it relieved symptoms of irritable bowel syndrome with fewer side effects than available drugs (T. J. Kaptchuk *et al. Br. Med. J.* **336**, 999–1007; 2008). It is true that if a therapy cannot beat a fake version of itself in trials, it is not working as the therapist claims. But if it triggers a big enough placebo effect, it might still be the best treatment available.

If drugs are effective and placebo responses small, this does not matter much. But people tend to turn to alternative medicine for subjective, stress-related conditions such as chronic pain, depression, nausea and fatigue (all problems that can affect cancer patients in treatment). Drugs for these

conditions have significant downsides, such as unpleasant side effects and addiction, and placebo responses often account for most of the effect of the drug. So it becomes plausible that compared to popping a pill, a patient might get more relief — and fewer side effects — from an hour with a sympathetic therapist.

The benefits of therapies such as reiki and acupuncture go beyond what we normally think of as placebo effects, however. Alternative therapists do not get results just because they are particularly good at fooling people into thinking that they will get better. Many elements of the care they provide — from talking to touch — seem to have the power to relieve symptoms and even influence physical outcomes. These elements do not show up when therapies are compared against sham treatments, because they are present in both arms of a trial.

Such benefits can be indirect. For example, tackling patients’ anxiety during invasive procedures such as keyhole surgery can reduce the risk of dangerous fluctuations in heart rate. This results not only from the direct effects on physiology, but also probably from patients needing lower doses of sedatives and painkillers.

Conventional medicine, with its squeezed appointment times and overworked staff, often struggles to provide such human aspects of care. One answer is to hire alternative therapists.

This ensures that such therapies are regulated, and that patients also get the conventional treatment they need. Such ‘integrative medicine’ is now offered by dozens of major US academic medical institutes. The Stanford Center for Integrative Medicine in California offers acupuncture to help

with chemotherapy side effects. If this helps patients to complete a conventional treatment by making those symptoms bearable, one therapist there told me, it might improve survival rates, too.

Critics say that this is dangerous quackery. Endorsing therapies that incorporate unscientific principles such as auras and energy fields encourages magical thinking, they argue, and undermines faith in conventional drugs and vaccines. That is a legitimate concern, but dismissing alternative approaches is not evidence-based either, and leaves patients in need.

Instead of rejecting such approaches wholesale, let’s learn from them. That means going beyond the simplistic practice of jettisoning anything that cannot beat placebo. We must tease out the real active ingredients of these therapies — things such as ritual, mental imagery, empathy, care and hope — so that we can learn how they work and find ways to incorporate them into patient care. ■

Jo Marchant is a science journalist and author of *Cure: A Journey into the Science of Mind Over Body*, to be published in January 2016. e-mail: jomarchant26@yahoo.co.uk

CONVENTIONAL
MEDICINE
OFTEN
STRUGGLES
TO PROVIDE
HUMAN
ASPECTS OF CARE.

➔ **NATURE.COM**
Discuss this article
online at:
go.nature.com/xp26zy

RESEARCH HIGHLIGHTS

Selections from the
scientific literature

IMMUNOLOGY

Infections 'scar' immune system

After the body has cleared a gut bacterial infection, some intestinal tissues have long-lasting weakened immunity — partly because of gut microbes.

Infection can lead to chronic inflammatory disorders. To find possible mechanisms, Yasmine Belkaid at the National Institute of Allergy and Infectious Diseases in Bethesda, Maryland, and her colleagues infected mice with a foodborne pathogen, *Yersinia pseudotuberculosis*, and then monitored them for up to six months. They found a variety of changes that weaken the immune system in gut tissues. For example, immune cells called dendritic cells were diverted into fat tissue away from lymph nodes, where they would normally activate the immune response. Signals from gut microbes also seemed to maintain inflammation.

The results show how frequent infections could lead to chronic disease later in life. *Cell* 163, 354–366 (2015)

CHEMISTRY

Cheap absorber for solar cells

An iron-based chemical can absorb and convert light into electrons with 92% efficiency — making it a promising material for solar cells.

Light-harvesting 'sensitizers' in solar cells are typically made of rare elements, which are expensive to scale up. To find a cheaper alternative, Kenneth Wärnmark at Lund University in Sweden and his colleagues engineered an iron-based sensitizer that captures photons and transfers their energy to electrons in a similar way to those based on rare elements.



METEOROLOGY

Winged weather watchers

Soaring birds equipped with sensors that track their location could be used to estimate weather variables such as wind speed.

Jelle Treep at the University of Amsterdam and his team analysed Global Positioning System (GPS) data from four griffon vultures (*Gyps fulvus*; pictured) as the birds flew across the Grand Causses region of southern France. By tracking the birds' in-flight location at three-second intervals and using knowledge of airborne bird movements, the team estimated wind speed and direction and vertical air movement. These GPS estimates agreed with ground data at three local weather stations and were consistent with meteorological theories.

As GPS equipment becomes smaller, lighter and able to collect higher-resolution data, bird-borne trackers could become an important tool for meteorological surveys in remote areas, the authors say.

Bull. Am. Meteorol. Soc. <http://doi.org/768> (2015)

The sensitizer converts photons with 92% efficiency: 12% better than the previous best sensitizer based on iron.

Some electrons quickly combined with positive charges, limiting the effective current generated. Still, the authors say that using abundant materials such as iron as a sensitizer for photovoltaics opens up possibilities for

low-cost solar energy.

Nature Chem. <http://dx.doi.org/10.1038/nchem.2365> (2015)

PALAEOGENETICS

First ancient African genome

A 4,500-year-old human skeleton from a cave in Ethiopia has produced Africa's

first ancient genome sequence.

Marcos Gallego Llorente at the University of Cambridge, UK, and his colleagues sequenced genomic DNA from one of the bones and compared it with sequences from modern Africans and Eurasians, as well as ancient Europeans. They conclude that the ancestors of modern Ethiopian highlanders were related to early farmers who moved into Europe from western Eurasia around 9,000 years ago. Descendants of these people later moved back to Africa around 3,000 years ago.

Most Africans today have 4–7% Eurasian ancestry because of this migration, suggesting that this event was larger and more significant than was thought.

Science <http://doi.org/78d> (2015)

NEUROSCIENCE

People identified from brain activity

A map of connections between brain regions that are active during mental activity can be used as a unique, reproducible 'fingerprint' to identify individuals.

Emily Finn of Yale University in New Haven, Connecticut, and her colleagues studied data from the Human Connectome Project, which is mapping all of the structural and functional connections in the brain. They looked at 126 people whose brains were scanned while they were resting or doing certain tasks. By analysing patterns of neural connectivity, the team identified subjects with a success rate of more than 90% when comparing rest scans, and with 54–87% success when comparing brain activity during tasks. The most useful networks for identifying people were those in certain regions of the cerebral cortex that control attention, memory and other

CHRISTIAN AUSSAGUEL

cognitive functions.

The results provide a foundation for future work to link functional brain connections with individual behaviours, the authors say. *Nature Neurosci.* <http://dx.doi.org/10.1038/nn.4135> (2015)

PALAEOANTHROPOLOGY

Early human with a familiar handshake

A recently discovered early human species probably walked upright and wielded tools, but also took to the trees.

Last month, researchers reported the discovery of fossil bones from at least 15 individuals of a species they named *Homo naledi*. A team led by Tracy Kivell at the University of Kent in Canterbury, UK, has analysed nearly 150 hand bones from the find, including a complete right hand (pictured, left). The hands resemble those of *Homo sapiens*, Neanderthals and other regular tool-users, although the long, curved fingers suggest that *H. naledi* was comfortable in trees.

In a separate study, William Harcourt Smith at the City University of New York and Jeremy DeSilva at Dartmouth College in Hanover, New Hampshire, looked at 107 foot bones, including a nearly complete right foot (pictured, right), and concluded that *H. naledi* strode upright. However, the feet still had some primitive features: certain toe bones were more curved than are those of modern humans.

Nature Commun. 6, 8431; 6, 8432 (2015)



BLUE BRAIN PROJECT

PETER SCHMIDT/
WILL HARCOURT-SMITH

CANCER

How elephants dodge cancer

Elephants have extra copies of a gene that fights tumour cells, which could explain why they rarely develop cancer.

Joshua Schiffman at the University of Utah in Salt Lake City and his colleagues studied elephant white blood cells and found that they have 20 copies of a tumour-suppressor gene called *TP53* in their genome — humans and other mammals have only one. The cells also underwent *TP53*-mediated apoptosis — programmed cell death — more frequently than human cells do when exposed to DNA-damaging radiation. This suggests that elephant cells kill themselves to avoid the risk of uncontrolled growth.

In a separate study, Vincent Lynch at the University of Chicago in Illinois and his co-workers report similar results. They also discovered more than a dozen *TP53* copies in two extinct species of mammoth, but just one copy in manatees and in small furry mammals called hyraxes — both close living relatives of elephants. The extra copies may have evolved as the animals in the elephant lineage expanded in size, says the team.

J. Am. Med. Assoc. <http://doi.org/772> (2015); Preprint at [bioRxiv](http://bioRxiv.org/773) <http://doi.org/773> (2015)

CLIMATE-CHANGE BIOLOGY

Corals cope with pH-altered waters

Some corals seem to be resilient to ocean acidification.

As carbon dioxide emissions rise, ocean waters are absorbing more of the gas and becoming less alkaline, threatening the ability of corals and other marine organisms to make skeletons and shells. Lucy Georgiou at the University of Western Australia in Perth and her colleagues exposed colonies of *Porites cylindrica* coral on Australia's Great Barrier Reef to flumes of modified

SOCIAL SELECTION

Popular topics
on social media

Nobel prizes prompt surprise online

This year's Nobel prizewinners seemed to surprise many researchers, judging by their reactions on social media. The chemistry prize recognized discoveries in DNA repair, yet was not awarded to the scientists who won the prestigious Lasker prize earlier this year for research in a similar area. Many speculated that the physics prize would go to a woman for the first time in more than 50 years, but it went to two men for their work on neutrinos instead. One of the scientists who shared the Nobel Prize in Physiology or Medicine, Chinese pharmacologist Youyou Tu, discovered a malaria medicine called artemisinin after studying traditional Chinese medicine texts. As chemist Ashutosh Jogalekar wrote on his blog (go.nature.com/loead): "The story of artemisinin clearly indicates that we need to pay much more attention to forgotten examples from traditional Asian medicine and subject them to scrutiny."

➔ **NATURE.COM**
For more on
popular papers:
go.nature.com/d2hnyf

sea water. This lowered the ambient pH around the animals so that it was similar to conditions that are predicted for oceans at the end of the century. After six months, the researchers found no difference in the growth rate of the corals' skeletons between controls and those living in lower pH conditions.

The corals naturally produce a fluid that bathes the growing parts of their skeletons, and the team found that the fluid had a higher pH than the reef waters for all the corals in the experiment. This suggests that some corals can regulate their internal pH to tolerate a certain level of ocean acidification, the authors say.

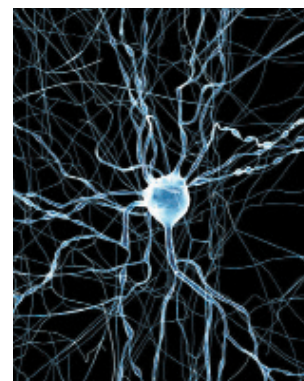
Proc. Natl Acad. Sci. USA <http://doi.org/77b> (2015)

NEUROSCIENCE

Computer model of rat-brain part

A supercomputer has simulated the activity of 31,000 virtual brain cells (pictured) in a section of rat brain the size of a grain of sand.

Henry Markram at the Swiss Federal Institute of Technology in Lausanne and his team built their model based on experimental measurements



of rat brain slices. The simulation represents roughly 37 million synapses, or neuronal connections, in the brain region that receives sensory information from the whiskers and other parts of the body. Using the model, the team simulated rat whisker movement and saw similar neuronal responses to those observed in rat experiments.

The model could be manipulated in ways that are difficult to do experimentally, providing insights into how individual cells contribute to the functions of neuronal networks, the authors say. *Cell* 163, 456–492 (2015)

➔ **NATURE.COM**
For the latest research published by
Nature visit:
www.nature.com/latestresearch

CORRECTION

The print version of the Research Highlight 'Corals cope with acidified waters' (*Nature* **526**, 296–297; 2015) incorrectly stated that ocean water is being acidified when in fact it is becoming less alkaline; the online title was changed to reflect that. It also said coral-made fluid was less acidic than reef waters; in fact, the fluid had a higher pH. And it said that some corals can control the pH of surroundings, whereas they control their internal pH.

SEVEN DAYS

The news in brief

EVENTS

Gene patent invalid

The High Court of Australia has ruled a key gene patent invalid, casting doubt on the status of other gene patents worldwide. On 7 October, the court determined that isolated DNA corresponding to a cancer-associated gene called *BRCA1* was not “a manner of manufacture”, and thus could not be patented under Australian law. The decision puts an end to a five-year legal battle. The US Supreme Court ruled a similar patent invalid in 2013.

Start-up lab

The venture-capital firm Y Combinator (YC) in Mountain View, California, is establishing a non-profit arm devoted to basic research in fields ranging from information technology to biology. YC Research will employ its own scientists and will initially consist of just one lab. YC president Sam Altman says that he is contributing US\$10 million of his own money to the enterprise, but that most details are still confidential.

Orca breeding ban

The California Coastal Commission has banned the breeding of killer whales (*Orcinus orca*) at the SeaWorld theme park in San Diego. SeaWorld was granted permission to expand its orca tanks in the park, but said that it was “disappointed” by the commission’s conditions, which company president Joel Manby said would be “inhumane” in depriving the animals of “the natural and fundamental right to reproduce”. SeaWorld has been increasingly targeted by campaigners who say that keeping killer whales in captivity is cruel; the company



XL CATLIN SEAVIEW SURVEY

Sea warming triggers coral bleaching

The US National Oceanographic and Atmospheric Administration (NOAA) stated on 8 October that Earth’s oceans are experiencing a mass coral-reef bleaching (pictured here in fire coral in Bermuda). Sea surface temperatures have risen enough to bleach reefs in three major coral-containing basins — the Pacific,

Indian and Atlantic oceans. It is the third global bleaching event in recorded history, after those in 1998 and 2010; as in the previous events, the El Niño weather pattern is helping to keep ocean temperatures high. By the end of 2015, 38% of the world’s coral reefs could be affected, NOAA says. See go.nature.com/1h3bmX for more.

denies these claims. The ruling does not affect SeaWorld sites outside California.

system, which explains how cells use enzymes to repair damage caused by ultraviolet light. See page 307 for more.

Washington DC to undertake the survey. See go.nature.com/ddo2bz for more.

AWARDS

DNA-repair Nobel

The 2015 Nobel Prize in Chemistry was awarded to three scientists who mapped how cells repair damaged DNA. Tomas Lindahl, Paul Modrich and Aziz Sancar shared the prize, announced on 7 October. Each discovered a different molecular process. Lindahl described how enzymes seek, cut out and patch up sections of damaged DNA, a mechanism called base-excision repair. Modrich worked on mismatch repair, which sorts out errors that are introduced when DNA is copied. And Sancar contributed research on the nucleotide-excision repair

RESEARCH

Cell lines unchecked

More than half of biomedical researchers do not authenticate the cell lines that they use in their experiments, a survey published on 12 October shows (L. P. Freedman *et al.* *BioTechniques* **59**, 189–192; 2015). Of the 446 scientists surveyed, 52% said that they do not do checks for species, tissue type and sex, with many citing cost and time constraints as the reason. Concern about reproducibility of results and wasted research funds because of misidentified cells led the non-profit Global Biological Standards Institute in

FACILITIES

Telescope start

Construction began on 9 October of a 23-metre telescope at Roque de los Muchachos Observatory on La Palma, in Spain’s Canary Islands. The instrument is destined to be part of the world’s largest γ -ray observatory, the Cherenkov Telescope Array, an international project expected to cost more than US\$300 million, with its Southern Hemisphere counterpart in Paranal, Chile. The dish will detect faint flashes of Cherenkov radiation — blue light emitted by the showers of electrons that

γ -ray photons unleash when they hit the atmosphere. The telescope is expected to start operating by late 2017.

Community lasers

A powerful laser facility near Prague is due to open officially on 19 October. The €250-million (US\$284-million) centre in the Czech Republic marks the first installation of the three-pronged Extreme Light Infrastructure (ELI), an experiment by the European Union to build large research facilities in countries that could not ordinarily afford them (see *Nature* **500**, 264–265; 2013). The ELI will largely be paid for by European structural funds, which typically finance civic projects such as road repair and waste clean-up. It is set to include two other laser facilities in Hungary and Romania.

PEOPLE

Richard Heck dies

Nobel-prizewinning chemist Richard Heck (**pictured**) died on 9 October, aged 84. Heck was a professor emeritus at the University of Delaware in Newark, and lived in Manila. He shared the 2010 Nobel Prize in Chemistry with Ei-ichi Negishi and Akira Suzuki for work he did in the 1960s and 1970s on reactions that link two carbon atoms together using a palladium catalyst. His palladium cross-coupling



reaction, known as the Heck reaction, opened up a way for chemists in fields from pharmaceuticals to electronics to make myriad molecules more easily.

Sexual harassment

Geoff Marcy, an astronomer at the University of California, Berkeley, renowned for his work on exoplanets, has been found to have violated campus sexual-harassment policy. BuzzFeed News revealed the findings, which were confirmed to *Nature* by a university spokesperson, on 9 October. The investigation was triggered by formal complaints, and found that Marcy had violated policies in a number of incidents involving students between 2001 and 2010. Marcy will face consequences — “sanctions that could include suspension or dismissal”, according to the university — only if he continues to harass students.

Marcy apologized in an open letter to the astronomy community. On 12 October members of the university's astronomy department released a statement advocating that Marcy be removed from the faculty. See go.nature.com/wqv2ng for more.

BUSINESS

CRISPR alliance

Genome-editing company Caribou Biosciences of Berkeley, California, and industrial chemicals giant DuPont of Wilmington, Delaware, announced on 8 October that they would join forces to exploit CRISPR/Cas9, a system used to make targeted changes to genomes. DuPont will gain exclusive rights to use Caribou's patented CRISPR/Cas9 technology in certain crops, and plans to bring genome-edited crops to market within the next decade.

Heart drug halted

Drug company Eli Lilly announced on 12 October that it was prematurely halting the phase III trial of its heart-disease drug evacetrapib. The drug, developed to treat people with a high risk of their arteries hardening and narrowing, is part of a family of cholesteryl ester transfer protein inhibitors, which increase fat-removing particles called high-density

COMING UP

14–16 OCTOBER

The European Food Safety Authority holds its second scientific meeting as part of the World Expo 2015 in Milan, Italy.

go.nature.com/zbihtw

14–17 OCTOBER

Dallas, Texas, hosts the Society of Vertebrate Paleontology's 75th anniversary meeting.

go.nature.com/hmmzkm

20–22 OCTOBER

The 11th World Conference on Bioethics, Medical Ethics and Medical Law takes place in Turin, Italy.

go.nature.com/ammaum

lipoproteins. The company stopped the trial after a review found that evacetrapib was ineffective; there were no safety concerns about the drug. The news triggered an 8% drop in Lilly's share price on the day of the announcement.

Toxic water

A US federal jury has found industrial chemicals giant DuPont in Wilmington, Delaware, liable for contaminating drinking water with C-8, an ingredient used in the manufacture of Teflon. Carla Bartlett, who developed kidney cancer after drinking the water, was awarded US\$1.6 million in damages. Her lawsuit — the first of some 3,500 similar ones to go to trial — claimed that the water caused the cancer. The 7 October ruling found no malice on the part of DuPont, which used the chemical at its plant in Parkersburg, West Virginia. A DuPont spin-off company, Wilmington-based Chemours, will bear the liability.

NATURE.COM

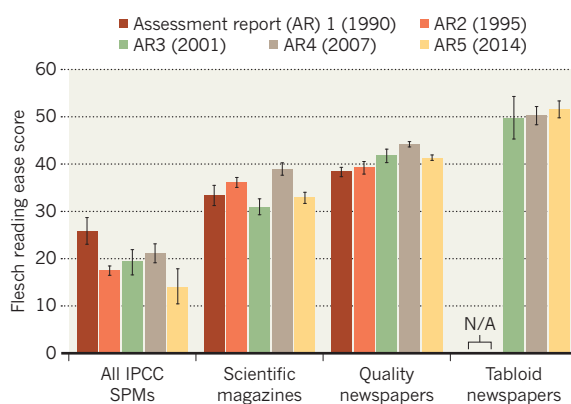
For daily news updates see:
www.nature.com/news

TREND WATCH

The summary findings of the Intergovernmental Panel on Climate Change (IPCC) are becoming increasingly unreadable, a study shows (R. Barkemeyer *et al.* *Nature Clim. Change* <http://doi.org/79f>; 2015). The IPCC's first assessment report, from 1990, scored highest in the Flesch reading ease test, which analyses word complexity and sentence structure. Its 2014 report scored lowest. By contrast, news stories about the IPCC are becoming easier to read. See go.nature.com/4ithkr for more.

COMPLEX CLIMATE REPORTS

The IPCC's summaries for policymakers (SPMs) are getting harder to read, although media stories about the reports are written more clearly.



NEWS IN FOCUS

HOT GENOMES Ancient DNA sequencing becomes less Eurocentric **p.303**

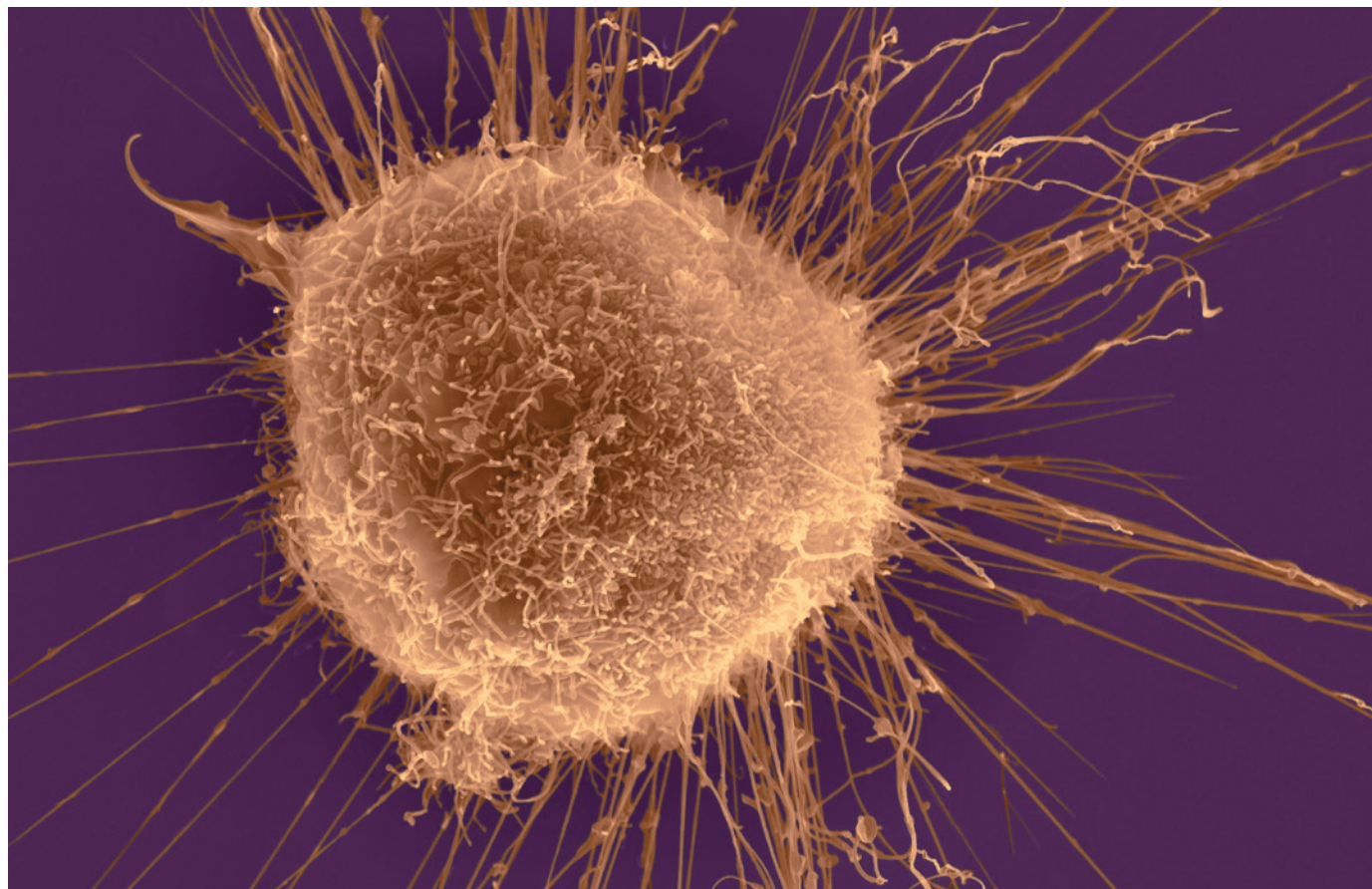
METROLOGY Experiments to redefine kilogram agree at last **p.305**

CLIMATE Companies that suck carbon from air get customers **p.306**



NEUROSCIENCE Do gut microbes shape brain development? **p.312**

DENNIS KUNKEL MICROSCOPY/VISUALS UNLIMITED/CORBIS



The challenges include mapping the diversity of cells in a tumour, such as this breast-cancer cell.

FUNDING

Giant charity lays out grand challenges of cancer

Cancer Research UK invests £100 million in fashionable way to distribute cash.

BY HEIDI LEDFORD

How can you distinguish between deadly cancers that need treatment and benign cancers that do not? Can unknown carcinogens be identified by the molecular fingerprints that they leave in tumours? Could vaccines be used to prevent more cancers?

The world's largest cancer charity hopes to find out. On 12 October, Cancer Research UK (CRUK) in London announced its intention to dedicate at least £100 million (US\$153 million) over 5 years to funding research teams to take on these and four other 'grand challenges'. The targets were identified by a panel of independent international researchers

convened by the charity, as well as input from the wider community. Teams eligible for funding can have members from anywhere in the world, but must include a UK-based contingent.

The charity will continue to fund projects proposed by investigators. But its new grand-challenge approach is in keeping with the ►

GRAND CHALLENGES

Things to know about cancer

Cancer Research UK intends to spend £100 million (US\$153 million) in pursuit of seven goals that could spur advances in preventing and treating the disease. They are:

- **Develop vaccines to prevent non-viral cancers.** This could stop a tumour early, before it becomes well established, metastasizes or begins to suppress the immune system.
- **Eradicate EBV-induced cancers from the world.** Epstein-Barr virus causes 200,000 cancers per year. Vaccines or drugs that target this virus could save many lives.
- **Discover how unusual patterns of mutation are induced by different cancer-causing events.** Sometimes a carcinogen will produce a distinct pattern of DNA mutations inside cells. Such patterns could reveal previously unknown causes of cancer.
- **Distinguish between lethal cancers that need treating, and non-lethal cancers that don't.** Catching cancer early is important

in defeating the disease, but over-diagnosis leads to unnecessary treatments — some of which can be life-altering, or have dangerous side effects.

- **Find a way of mapping tumours at the molecular and cellular level.** Tumours are made of a genetically diverse mixture of cells. Understanding what cells are present and what they are doing could enable physicians to predict how a tumour will respond to a given therapy.
- **Develop innovative approaches to target the cancer super-controller MYC.** The protein encoded by the *MYC* gene is often considered undruggable, but new approaches could change that.
- **Deliver biologically active macromolecules to any and all cells in the body.** Large molecules such as proteins and RNA can be remarkably specific drugs — but getting them into cells is a challenge that has dogged the field. **H.L.**

► times, says Nic Jones, the charity's chief scientist. "There's a cultural change happening within biological sciences," he says. "The idea is to be bold and take risks, and do things that might not work and might take a long time."

That idea is sweeping the globe: other prominent efforts have issued 'grand challenges' that have tasked researchers with, among other things, devising ways to track the rise of antibiotic resistance and keeping tabs on all asteroids that pose a threat to human populations.

The movement has its roots in a series of 23 challenges issued by mathematician David Hilbert in 1900; these are credited with setting the course for twentieth-century mathematics. The Bill & Melinda Gates Foundation in Seattle, Washington, resurrected the idea in 2003 with its vaunted Grand Challenges in Global Health programme, dedicating US\$450 million to achieve 14 broad goals that included creating new vaccines and improving nutrition. Since then, the foundation has awarded

1,752 grants in 81 countries under its grand-challenges banner.

The model has been taken up around the world by governments, charities and universities looking to tackle big societal problems — and get a little publicity in the process. "The term has become sexy," says Abdallah Daar, a public-health specialist at the University of Toronto in Canada who helped to craft the Gates Foundation challenges. "It has this element of bigness and difficulty and gets attention. And when people win, there's a big announcement."

WAYS FORWARD

Daar notes that the CRUK programme stands out for its focus on basic research. Although each of the challenges could eventually contribute to clinical innovations, most of them focus on laying the groundwork for therapeutic advances.

One challenge, for example, aims to develop ways to map the cellular make-up of a tumour

(see 'Things to know about cancer'). DNA sequencing studies have shown that the population of cells in a tumour is genetically diverse — and the genes lurking in those cells could determine whether the cancer will become resistant to specific drugs. Understanding the cellular topography could help physicians select the right treatment.

Another grand challenge asks researchers to find a way to deliver large molecules — such as certain proteins or RNAs — to any cell in the body. Although some large molecules, for example antibody therapies, are already used to treat cancer, they tend to operate outside the cell. "We know how to introduce these molecules in cell culture," says Tyler Jacks, director of the Koch Institute for Integrative Cancer Research at the Massachusetts Institute of Technology in Cambridge. "If we could only do that in the body, we might be able to open up the door to all sorts of new types of therapy."

THE UNDRUGGABLE

Other challenges are more specific: one, to eradicate cancers linked to the common Epstein-Barr virus, casts a spotlight on a cause of cancer that is particularly important in developing nations. And the charity's team of advisers has challenged researchers to design a drug that will block a cancer-associated protein called MYC. That is no simple task: researchers have been trying to do it for decades without success, and MYC is typically classified as an 'undruggable' target. But the hope, says Jones, is that methods uncovered in the attempt could also be applied to other elusive targets.

CRUK intends to issue new challenges each year, and Jones hopes that suggestions for the next round will come from the community. Crafting a tractable yet inspiring challenge is an art, says Richard Klausner, chief medical officer of biotechnology company Illumina in San Diego, California, and chair of the CRUK challenge-selection committee. A few of Hilbert's mathematical challenges remain unsolved more than a century after they were proposed; to avoid a similar outcome, Klausner says, the committee tried to avoid setting unattainable goals.

"You want it to be ambitious enough, but not crazy," he says. "They need to be at the edge of doable — but still doable." ■



**MORE
ONLINE**

TOP STORY



Scientists hope to attract millions to 'DNA.LAND' go.nature.com/pfncq2

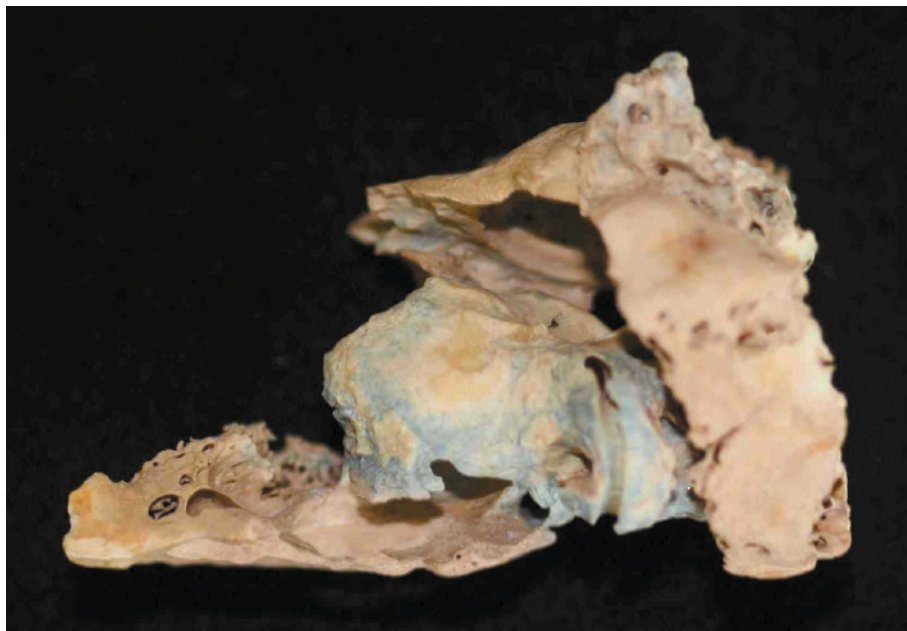
MORE NEWS

- Epigenetic 'tags' linked to homosexuality in men go.nature.com/kap88w
- How elephants avoid cancer go.nature.com/lqh3ok
- Fragment of rat brain simulated in supercomputer go.nature.com/cocjtg

NATURE PODCAST



Ancient human teeth in China, cooperative climate talks, and a worm surprises scientists nature.com/nature/podcast



The inner ear's petrous bone is a rich source of DNA in archaeological specimens.

PALAEoANTHROPOLOGY

Hot climes yield ancient DNA

After years of frustration, researchers are getting genetic material from old bones in warm places.

BY EWEN CALLAWAY

The quest to glean DNA from ancient humans started in hot climes. In 1985, a paper in *Nature* (S. Pääbo *Nature* **314**, 644–645; 1985) reported DNA sequences from an ancient Egyptian mummy — the first time anyone claimed to have isolated genetic material from a long-dead human.

But geneticist Svante Pääbo at the Max Planck Institute for Evolutionary Anthropology in Leipzig, Germany, later realized that the sequences were likely to be the result of contamination by modern DNA — possibly his own. Knowing that heat speeds up the decay of DNA, he and others turned their focus to remains from cooler climates, where DNA was more likely to persist.

And that meant, inevitably, a geographic bias in ancient DNA research. “It’s really been a bit Eurocentric,” says Hannes Schroeder, a palaeogeneticist at the Natural History Museum of Denmark in Copenhagen.

Now, after 30 years, researchers have returned to warmer places — this time finding more than contamination. Last week, a

team published Africa’s first ancient genome (M. Gallego Llorente *et al.* *Science* <http://doi.org/78d>; 2015), from a 4,500-year-old skeleton in Ethiopia, and the genetic profiles of dozens of 8,000-year-old skeletons from the Middle East were reported at a meeting. These results reveal lost migrations and colour in the details of murky prehistoric events far beyond Europe and Asia.

“There are so many more fascinating things happening in human world history than just in Europe,” says Schroeder, who is using ancient DNA to trace the origins of Caribbean slaves. Breakthroughs in the past decade have made it possible to study the very short and damaged DNA fragments typically found in remains from warmer places. Researchers have also found that a tiny part of the inner ear, the petrous portion of the temporal bone, holds much greater amounts of DNA than other human bones. “We find a crazy concentration of

ancient DNA in this region,” says Ron Pinhasi, an archaeologist at University College Dublin (R. Pinhasi *et al.* *PLoS ONE* **10**, e0129102; 2015).

Pinhasi and colleagues analysed DNA from the petrous of 26 individuals who lived in northwestern Anatolia (present-day Turkey) around 8,000 years ago, about the time that farming spread from the Middle East into Europe. The analysis suggests that they belonged to the source of a migration into southeast and central Europe that displaced local hunter-gatherers, the team reported on 9 October at a meeting of the American Society of Human Genetics in Baltimore, Maryland (I. Matheison *et al.* Preprint at bioRxiv at <http://doi.org/4wt>; 2015).

Pinhasi also turned to the petrous bone after getting access to a 4,500-year-old human skeleton excavated in 2012 in the Mota cave in the highlands of Ethiopia. As reported in *Science* last week, he and his colleagues recovered enough DNA to sequence the ancient Ethiopian’s genome 12 times over, producing the first complete genome from an ancient African.

The man’s genome is more closely related to present-day Ethiopian highlanders known as the Ari than to any other population the team examined. And using genetic evidence from Eurasian ancient genomes and present-day Eurasian and African populations, the researchers determined that the Ari descend from Middle Eastern farmers. They suggest that after expanding into Europe, some of these people later moved south to Africa, bringing new crops (see <http://dx.doi.org/79n>; 2015).

“Africa is the next step in terms of the population history of humans. It is an obvious place to go and look,” says Eske Willerslev, an evolutionary geneticist at the Natural History Museum of Denmark. He says that his team has also had success in obtaining DNA from African remains, including some that have not been protected by the relatively cool confines of a high-altitude cave. Willerslev expects that researchers will eventually strike lucky in extracting DNA from African remains that are many tens of thousands of years old, rather than a few thousand.

Some researchers speculate that remains around 100,000 years old from the Skhul and Qafzeh caves in Israel might represent failed migrations out of Africa that preceded the global expansion of *Homo sapiens* by tens of millennia, whereas others disagree. Ancient DNA may be the only way to settle the debate.

And few skeletons provoke such a debate among ancient-DNA researchers as *Homo floresiensis*, the mysterious ‘hobbit’ fossil discovered a decade ago on the Indonesian island of Flores. Past efforts to reap DNA from its bones have failed, yet its genetics may be the only way to determine its evolutionary relationship to *H. sapiens*. “I think there’s a shot with *floresiensis*,” Pinhasi says. “I think that has to be tested.” ■

Additional reporting by Sara Reardon.

METROLOGY

Experiments to redefine kilogram converge at last

After a fraught few years, results agree in time to meet a 2018 deadline.

BY ELIZABETH GIBNEY

For decades, metrologists have strived to retire 'Le Grand K' — the platinum and iridium cylinder that for 126 years has defined the kilogram from a high-security vault outside Paris. Now it looks as if they at last have the data needed to replace the cylinder with a definition based on mathematical constants.

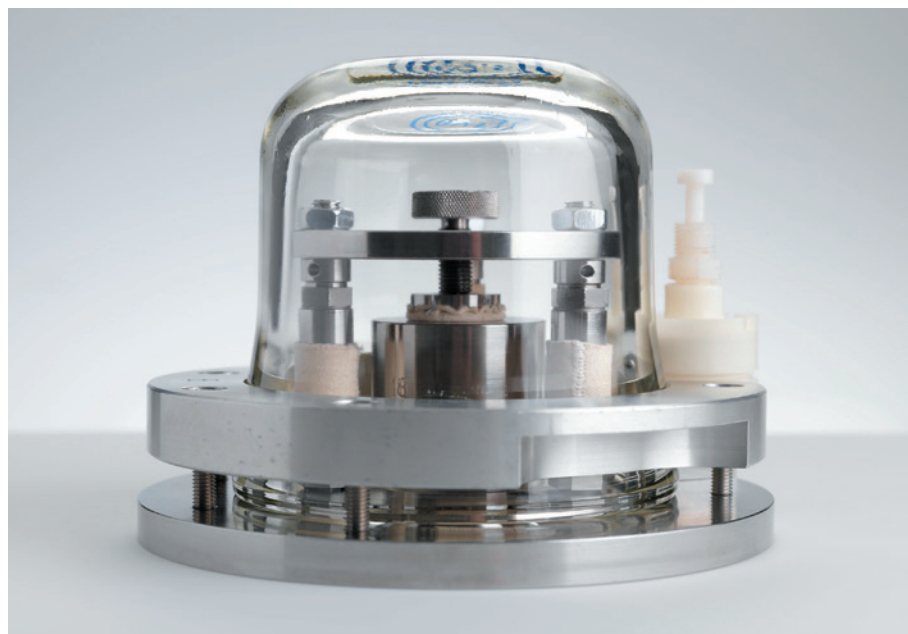
The breakthrough comes in time for the kilogram to be included in a broader redefinition of units — including the ampere, mole and kelvin — scheduled for 2018. And this week, the International Committee for Weights and Measures (CIPM) will meet in Paris to thrash out the next steps.

"It is an exciting time," says David Newell, a physicist at the US National Institute of Standards and Technology (NIST) in Gaithersburg, Maryland. "It is the culmination of intense, prolonged efforts worldwide."

The kilogram is the only SI unit still based on a physical object. Although experiments that could define it in terms of fundamental constants were described in the 1970s, only in the past year have teams using two completely different methods achieved results that are both precise enough, and in sufficient agreement, to topple the physical definition.

Redefinition will not make the kilogram more precise, but it will make it more stable. A physical object can lose or gain atoms over time, or be destroyed, but constants remain the same. And a definition based on constants would, at least in theory, allow the exact kilogram measure to be available to someone anywhere on the planet, rather than just those who can access the safe in France, says Richard Davis, former head of the mass division of the International Bureau of Weights and Measures (BIPM) in Sèvres, France, which hosts the metal kilogram.

In 2011, the CIPM formally agreed to express the kilogram in terms of Planck's constant, which relates a particle's energy to its frequency, and, through $E = mc^2$, to its mass. This means first setting the Planck value using experiments based on the current reference kilogram, and then using that value to define the kilogram. The CIPM's committee on mass recommends that three independent measurements of Planck's constant agree, and that two



A replica of the kilogram mass reference, which is set to be replaced by a definition based on constants.

of them use different methods.

One method, pioneered by an international team known as the Avogadro Project, involves counting the atoms in two silicon-28 spheres that each weigh the same as the reference kilogram. This allows them to calculate a value for Avogadro's constant, which the researchers convert into a value for Planck's constant.

"I think every metrologist worried, 'What if they never converge?'"

Another method uses a device called a watt balance to produce a value for Planck's constant by weighing a test mass calibrated according to the reference kilogram against an electromagnetic force. Reaching agreement proved difficult. In early 2011, some researchers contemplated simply averaging measurements from the two different devices (see *Nature* <http://doi.org/cjzwdn>; 2011). "I think every metrologist worried, 'What if they never converge?'" says Davis.

Such a fudge did not prove necessary, thanks to three years of intense work, says Joachim Ullrich, president of the German National Metrology Institute (PTB) in Braunschweig, which coordinates the Avogadro

Project, and chair of the CIPM's Consultative Committee for Units. The first sign of progress came after the Measurement Science and Standards laboratory in Ottawa, part of Canada's National Research Council (NRC), bought and rebuilt a watt balance originally constructed at the UK National Physical Laboratory in Teddington.

In a new lab, a fresh NRC team factored in some predicted but as yet unaccounted for systematic errors, and the result, published¹ in January 2012, inched closer to the Avogadro Project's silicon-sphere result.

That still left the result from NIST as an outlier, says Newell, chair of the international CODATA committee's group on fundamental constants, which provides a best value for constants such as Planck's every four years by taking into account the results of all the experiments done so far. "We brought in a whole new research team, we went over every component, went through every system," he says. They never found the cause for the disagreement, but in late 2014 the NIST team achieved² a match with the other two^{3,4}, who in the meantime had shrunk their relative uncertainties to within the required levels.

In August 2015, when CODATA ►

► published its latest value for Planck's constant, the uncertainty was 12 parts per billion, just over one-quarter of its value in CODATA's previous report — and within the CIPM's requirements.

The CIPM will discuss its next moves during its meeting at the BIPM on 15 and 16 October. This will include a discussion of the draft resolution that is expected to redefine the ampere, mole, kelvin and kilogram at the General Conference on Weights and Measures in 2018. The BIPM is still working on a protocol that will

allow teams without access to a watt balance or silicon-sphere set-up to use a new kilogram definition.

There is still scope for upset. The teams have until 1 July 2017 to publish further data before the value of Planck's constant is fixed. Before this deadline, Ullrich's team plans to use a new batch of spheres from Russia in experiments that he hopes will lead to even more-certain values for Planck's constant, but could cause the results to diverge again. "Then we would be in trouble," he says. "But I'm very confident this

will not happen." Newell agrees: "This train has a lot of momentum and there has to be something seriously wrong to derail it."

If they are proved right, in 2018, Le Grand K will join the metre as a museum piece. "We'll keep it," says Davis, "but it won't be defining anything anymore." ■

1. Steele, A. G. *et al. Metrologia* **49**, L8–L10 (2012).
2. Schlamminger, S. *et al. Metrologia* **52**, L5–L8 (2015).
3. Sanchez, C. A. *et al. Metrologia* **52**, L23 (2015).
4. Mana, G. *et al. J. Phys. Chem. Ref. Data* **44**, 031209 (2015).

CLIMATE CHANGE

Firms that suck carbon from air go commercial

Two companies announce that they are expanding and upgrading their plants.

BY DANIEL CRESSEY

It has long been regarded as one of the more blue-skies solutions to climate change. Now two companies have vastly increased their capability to suck carbon dioxide from the air. One, based in Canada, plans to convert captured CO₂ into diesel to fuel buses; the other, in Switzerland, will sell it on to a firm that uses CO₂ to boost crop growth in greenhouses.

The carbon emissions that this will save are not significant. But David Keith, executive chairman of the Canadian firm, Carbon Engineering in Calgary, and a climate physicist

at Harvard University in Cambridge, Massachusetts, says that his company's air-capture plant will position the technology to be further scaled up. Most significantly, he says, the plant will now run the whole process — from CO₂ capture to regeneration — for the first time.

Others are excited by the development. "The fact they're getting to commercial-scale prototypes is incredibly encouraging," says Noah Deich, executive director of the Center for Carbon Removal in Berkeley, California.

More than a dozen facilities worldwide, including oil refineries and power plants, already capture millions of tonnes of CO₂ from

the flue gases they expel. The idea of capturing carbon directly from the atmosphere — where CO₂ is present in much lower concentrations than in flue gases and so is harder to extract — has been around for several years, but only in the form of small, demonstration projects.

On 9 October, Carbon Engineering officially opened a new plant in Squamish, British Columbia, that can capture and process around 1 tonne of CO₂ per day — about the same as a typical car might emit when driven about 5,000 kilometres. This represents a big step up from the company's earlier demonstration plant, which ran only the first step of capture and did not regenerate gaseous CO₂.

The plant uses fans to push air through towers containing potassium hydroxide solution, which reacts with CO₂ to form potassium carbonate; the remaining air, now containing less CO₂, is released. Further treatment of the solution separates out the captured CO₂, regenerating the capture solution for reuse. These processes are currently powered by electricity, which in British Columbia is mainly generated by hydroelectric sources, says Keith. Initially, the company will re-release the captured CO₂, but Carbon Engineering announced last week that it had signed a Can\$435,000 (US\$333,000) deal with the province of British Columbia to assess the potential of turning the CO₂ into fuel to power local buses.

Meanwhile, the Swiss company, Climeworks in Zurich, announced at a UK meeting on greenhouse-gas capture in Oxford earlier this month that it plans to start capturing CO₂ on a commercial scale. Its plant in Hinwil, Switzerland, will capture 1,000 tonnes of CO₂ per year starting in mid-2016, according to Anca



Carbon Engineering's demonstration plant in British Columbia captures carbon dioxide from the air.

CARBON ENGINEERING

Timofte, a process engineer at the company.

In some ways the technology is similar to that of Carbon Engineering, but Climeworks will instead use granules to soak up the CO₂, using a module that will sit on top of an incineration plant. (The technology is still classed as air capture because the material will scrub CO₂ from air near the plant rather than from the expelled gases.) Waste heat from the incinerator will be used to drive the captured CO₂ off the granules, which can then be reused.

The company has arranged to sell CO₂ produced in this way to the firm Gebrüder Meier, which will use it to increase crop yields in greenhouses. Climeworks is also assessing the beverage industry as a source of potential customers, says Timofte.

If such companies are to scale up further and make money, one challenge will be finding buyers for their CO₂, says Tim Kruger, a geoengineer at the University of Oxford, UK, who organized the Oxford meeting and runs a company, Origen Power, that hopes to generate carbon-negative energy. And it is not clear whether companies will be able to produce CO₂ or related products at a price that is competitive enough to attract a wide pool of clients.

In 2011, a report from the American Physical Society (APS) estimated that air capture



David Keith, chairman of Carbon Engineering.

would cost at least US\$600 per tonne of CO₂, assuming a large system that removed 1 million tonnes of CO₂ per year. But Climeworks says that its price will be in that range in the first year of its plant's operation, despite being on a smaller scale than the APS example. The company also expects that cost to fall as the technology develops. Keith, meanwhile, says that the CO₂ produced by Carbon Engineering's plant is expensive, but emphasizes that it is a pilot; he says that prices of \$100–200

per tonne of CO₂ are realistic for the bigger iterations that it is planning.

Even if the companies cannot compete on price with conventionally manufactured CO₂ (which can be as low as tens of dollars per tonne but can be significantly higher), there are other factors that could help to create demand for air-captured CO₂. The introduction of a carbon tax could incentivize big emitters to pay other companies to mop up their CO₂ to avoid paying the tax. And if the world is ever to become completely carbon neutral, air capture will have a part to play, says Nilay Shah, an engineer working on low-carbon technologies at Imperial College London.

Efforts to mitigate climate change should focus on capturing CO₂ at the source. But there are many scenarios in which pre-emission capture is not viable. "Once you start to get into things like capturing carbon from vehicles or from household boilers, that's much more expensive," says Shah. "You may well be better off capturing CO₂ from the air."

Keith emphasizes that his company is not trying to fix climate change on its own. "Air capture has been stuck in a catfight between one group of people saying it's a silver bullet and one group saying it's bullshit," he says. "The truth is it's neither." ■

AWARDS

DNA-repair sleuths win chemistry Nobel

Tomas Lindahl, Paul Modrich and Aziz Sancar share prize for work on how DNA heals itself.

BY DANIEL CRESSEY

The 2015 Nobel Prize in Chemistry was awarded last week to three researchers for their work on DNA repair.

Tomas Lindahl, Paul Modrich and Aziz Sancar "mapped, at a molecular level, how cells repair damaged DNA and safeguard the genetic information", said the Royal Swedish Academy of Sciences in Stockholm, which awards the prize.

DNA is not a stable molecule, but slowly decays over time. For life to exist — as Lindahl first realized while working at the Karolinska Institute in Stockholm in the 1970s — there must be repair mechanisms that fight back against this process.

Numerous scientists have since chronicled the many ways in which damaged DNA is patched up, says Stephen West, who works on DNA repair at the Francis Crick Institute in London, where Lindahl is now an emeritus

group leader. "The DNA-repair field is a large field," says West. "Many of us thought a Nobel would not go to this field because there are so many people with a claim to the prize."

But the three repair mechanisms recognized with the Nobel prize "are probably the three most important and best-understood mechanisms", he says, adding that the awards are "fantastically well deserved".

REPAIR JOBS

Lindahl, who is regarded as one of the founders of the field, chronicled a process dubbed base excision repair, in which specific enzymes recognize, cut out and patch up bases in the DNA molecule. Before his work, "I don't think anybody really considered the idea that DNA requires active engagement by a set of house-keeping processes to keep it in a stable state," says Keith Caldecott, who studies DNA repair at the University of Sussex in Brighton, UK, and did postdoctoral work with Lindahl.

Sancar — who was born in Savur, Turkey, but has spent most of his professional life in the United States and is now at the University of North Carolina at Chapel Hill — worked in the 1980s to explain how cells use enzymes to repair damage to DNA from ultraviolet rays or other

"We need DNA repair but we don't like it that the cancer cells have DNA repair."

carcinogens, through a system called nucleotide excision repair. And in 1989, Modrich, who is at Duke University School of Medicine in Durham, North Carolina, published work on a third mechanism — 'mismatch repair' — which deals with errors produced when DNA is copied.

This September, the prestigious Albert Lasker Basic Medical Research Award was also awarded for work on how cells correct damage to DNA. But it went to two other researchers: Evelyn Witkin of Rutgers University in New ►



Tomas Lindahl, Paul Modrich and Aziz Sancar share the 2015 Nobel Prize in Chemistry.

Brunswick, New Jersey, and Stephen Elledge of Brigham and Women's Hospital in Boston, Massachusetts.

WIDER IMPACTS

Speaking to reporters in Stockholm at the Nobel press conference, Lindahl noted that understanding DNA repair has implications for human health. People with faults in their repair system have an increased risk of developing cancers, because damaging mutations

can go uncorrected. Cancer cells themselves survive damage by using enzymes to patch up DNA, and there is now interest in therapies that target DNA-repair pathways in tumour cells. "We need DNA repair but we don't like it that the cancer cells have DNA repair," Lindahl said.

Work in the field has had an impact in other areas, too. Lindahl's research proved influential in the 1980s and 1990s, when scientists were first working to extract and analyse ancient DNA. The patterns of DNA damage that he

first characterized are now used as a stamp of authenticity to show that DNA is ancient, and not modern contamination.

The chemistry prize follows the award of the Nobel Prize in Physiology or Medicine to William Campbell, Satoshi Ōmura and Youyou Tu for their work on therapies against parasitic infections. The physics Nobel went to Takaaki Kajita and Arthur McDonald for showing that neutrinos have mass. ■

CORRECTIONS

The News story 'Neutrino flip wins physics prize' (*Nature* **526**, 175; 2015) wrongly implied that physicists knew about all three types of neutrino in the 1960s. In fact, the tau neutrino was postulated only in the 1970s. The News Feature 'The impenetrable proof' (*Nature* **526**, 178–181; 2015) wrongly located the University of Antwerp in the Netherlands. It is, of course, in Belgium. In the News Feature 'The mitochondria mystery' (*Nature* **525**, 444–446; 2015), the quote "The standards for a shampoo seem to be harsher" was erroneously attributed to Ted Morrow. When he said these words, he was characterizing the stance of another researcher, not expressing his own opinion, and the quote should not have been included.

L TO R: JUSTIN TALLIS/AFP/GETTY; HHMI; SAMUEL CORUM/ANADOLU AGENCY/GETTY

THE LANDSCAPE FOR HUMAN GENOME EDITING

BY HEIDI LEDFORD

A view of international regulations suggests where in the world a CRISPR baby could be born.

They are meeting in China; they are meeting in the United Kingdom; and they met in the United States last week. Around the world, scientists are gathering to discuss the promise and perils of editing the genome of a human embryo. Should it be allowed — and if so, under what circumstances?

The meetings have been prompted by an explosion of interest in the powerful technology known as CRISPR/Cas9, which has brought unprecedented ease and precision to genetic engineering. This tool, and others like it, could be used to manipulate the DNA of embryos in a dish to learn about the earliest stages of human development. In theory, genome editing could also be used to ‘fix’ the mutations responsible for heritable human diseases. If done in embryos, this could prevent such diseases from being passed on.

The prospects have prompted widespread concern and discussion among scientists, ethicists and patients. Fears loom that if genome editing becomes acceptable in the clinic to stave off disease, it will inevitably come to be used to introduce, enhance or eliminate traits for non-medical reasons. Ethicists are concerned that unequal access to such technologies could lead to genetic classism. And targeted changes to a person's genome would be passed on for generations, through the germ line (sperm and eggs), fuelling fears that embryo editing could have lasting, unintended consequences.

Adding to these concerns, the regulations in many countries have not kept pace with the science.

Nature has tried to capture a snapshot of the legal landscape by querying experts and government agencies in 12 countries with histories of well-funded biological research. The responses reveal a wide range of approaches. In some countries, experimenting with human embryos at all would be a criminal offence, whereas in others, almost anything would be permissible.

Concerns over the manipulation of human embryos are nothing new. Rosario Isasi, a legal scholar at McGill University in Montreal, Canada, points to two key waves of legislation over the years: one sparked by concerns about the derivation of embryonic stem cells, which was largely deemed acceptable; the other about reproductive cloning, which was largely prohibited for safety reasons.

The current regulatory mosaic is their legacy. Tetsuya Ishii, a

bioethicist at Hokkaido University in Sapporo, Japan, spent nearly a year analysing relevant legislation and guidelines in 39 countries, and found that 29 have rules that could be interpreted as restricting genome editing for clinical use (M. Araki and T. Ishii *Reprod. Biol. Endocrinol.* 12, 108; 2014). But the ‘bans’ in several of these countries — including Japan, China and India — are not legally binding.

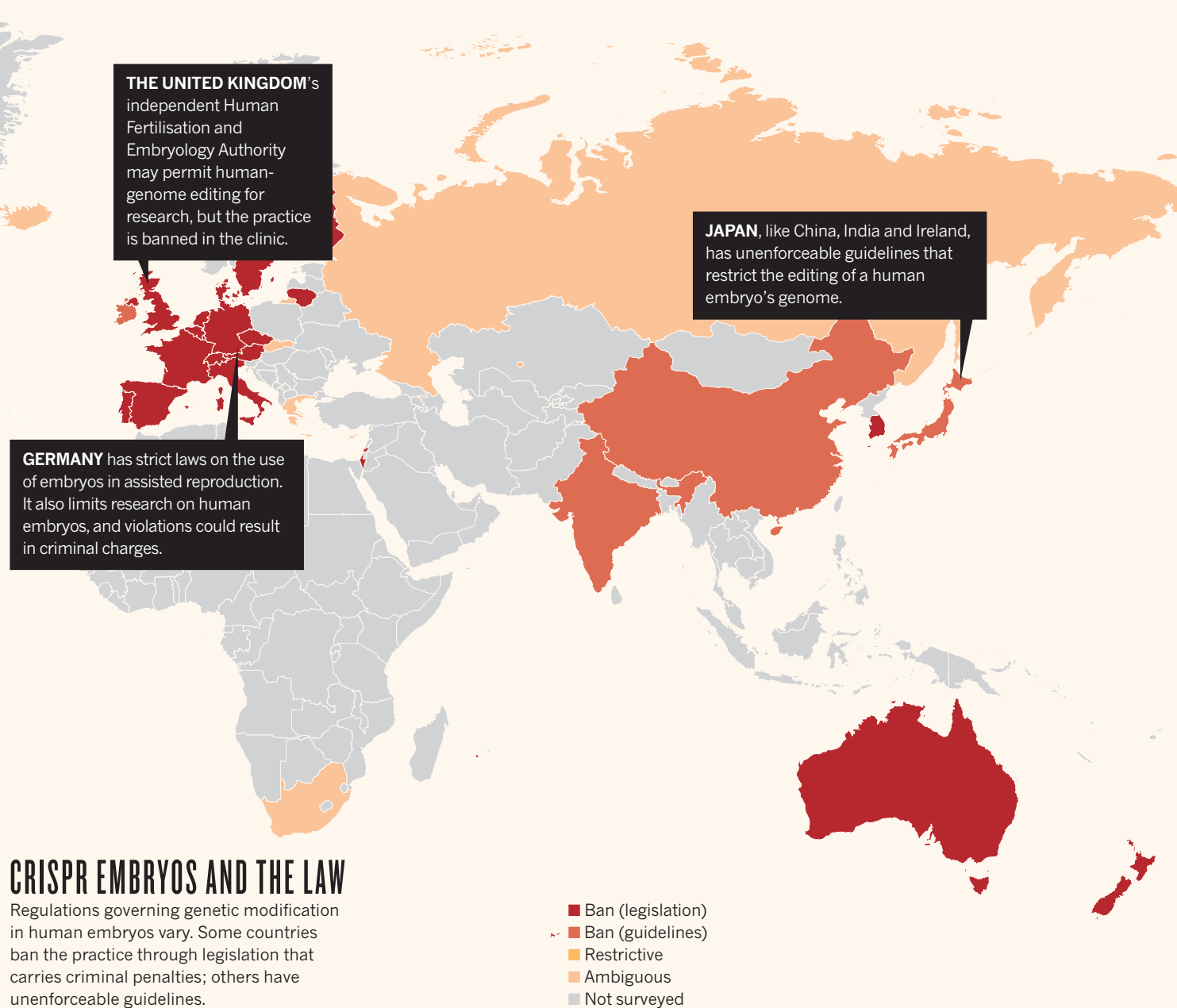
“The truth is, we have guidelines but some people never follow them,” said Qi Zhou, a developmental biologist at the Chinese Academy of Sciences Institute of Zoology in Beijing, at a meeting hosted by the US National Academy of Sciences in Washington DC last week. Ishii considers the rules in nine other countries — among them Russia and Argentina — to be “ambiguous”. The United States, he notes, prohibits federal funding for research involving human embryos, and would probably require regulatory approval for human gene editing, but does not officially ban the use of the technique in the clinic. In countries where clinical use is banned, such as France and Australia, research is usually allowed as long as it meets certain restrictions and does not

THE UNITED STATES does not allow the use of federal funds to modify human embryos, but there are no outright genome-editing bans. Clinical development may require approval.

ARGENTINA bans reproductive cloning, but research applications of human-genome editing are not clearly regulated.

**“WE HAVE GUIDELINES
BUT SOME PEOPLE NEVER
FOLLOW THEM.”**

SOURCE: M. ARAKI & T. ISHII *REPROD. BIOL. ENDOCRINOL.* 12, 108 (2014)



attempt to generate a live birth (see 'CRISPR embryos and the law').

Many researchers long for international guidelines that, even if not enforceable, could guide national lawmakers. Developing such a framework is one of the aims of ongoing discussions; the US National Academy, for example, plans to hold an international summit in December and then produce recommendations for responsible use of the technique in 2016.

But the research has already begun, and more is coming. Scientists in China announced in April that they had used CRISPR to alter the genomes of human embryos, albeit ones incapable of producing a live baby (P. Liang *et al.* *Protein Cell* 6, 363–372; 2015). Xiao-Jiang Li, a neuroscientist at Emory University in Atlanta, Georgia, who has used the technique in monkeys, says he has heard rumours that several other Chinese laboratories are already doing such experiments. And in September, developmental biologist Kathy Niakan of the Francis Crick Institute in London applied to the UK Human Fertilisation and Embryology Authority for permission to use the technique to study errors in embryo development that can contribute to infertility

and miscarriage. No one so far has declared an interest in producing live babies with edited genomes, and initial experiments would suggest that it is not yet safe. But some suspect that it is only a matter of time.

Ishii predicts that countries with high rates of *in vitro* fertilization will be the first to attempt clinical applications. Japan, he says, has one of the highest numbers of fertility clinics in the world, and has no enforceable rules on germline modification. The same is true for India.

Guoping Feng, a neuroscientist at the Massachusetts Institute of Technology in Cambridge, hopes that with improvement, the technique could eventually be used to prevent genetic disease. But he argues that it is much too soon to be trying it in the clinic. "Now is not the time to do human-embryo manipulation," he says. "If we do the wrong thing, we can send the wrong message to the public — and then the public will not support scientific research anymore." ■

**"IF WE DO THE WRONG
THING, WE CAN SEND THE
WRONG MESSAGE."**

Heidi Ledford writes for Nature from Cambridge, Massachusetts.

BRAIN, MEET GUT

BY PETER ANDREY SMITH

**Neuroscientists
are probing the
connections
between intestinal
microbes and brain
development.**

Nearly a year has passed since Rebecca Knickmeyer first met the participants in her latest study on brain development. Knickmeyer, a neuroscientist at the University of North Carolina School of Medicine in Chapel Hill, expects to see how 30 newborns have grown into crawling, inquisitive one-year-olds, using a battery of behavioural and temperament tests. In one test, a child's mother might disappear from the testing suite and then reappear with a stranger. Another ratchets up the weirdness with some Halloween masks. Then, if all goes well, the kids should nap peacefully as a noisy magnetic resonance imaging machine scans their brains.

"We try to be prepared for everything," Knickmeyer says. "We know exactly what to do if kids make a break for the door."

Knickmeyer is excited to see something else from the children — their faecal microbiota, the array of bacteria, viruses and other microbes that inhabit their guts. Her project (affectionately known as 'the poop study') is part of a small but growing effort by neuroscientists to see whether the microbes that colonize the gut in infancy can alter brain development.

The project comes at a crucial juncture. A growing body of data, mostly from animals raised in sterile, germ-free conditions, shows that microbes in the gut influence behaviour and can alter brain physiology and neurochemistry.

In humans, the data are more limited. Researchers have drawn links between gastrointestinal pathology and psychiatric neurological conditions such as anxiety, depression, autism, schizophrenia and neurodegenerative disorders — but they are just links.

"In general, the problem of causality in microbiome studies is substantial," says Rob Knight, a microbiologist at the University of California, San Diego. "It's very difficult to tell if microbial differences you see associated with diseases are causes or consequences." There are many outstanding questions. Clues about the mechanisms by which gut bacteria might interact with the brain are starting to emerge, but no one knows how important these processes are in human development and health.

That has not prevented some companies in

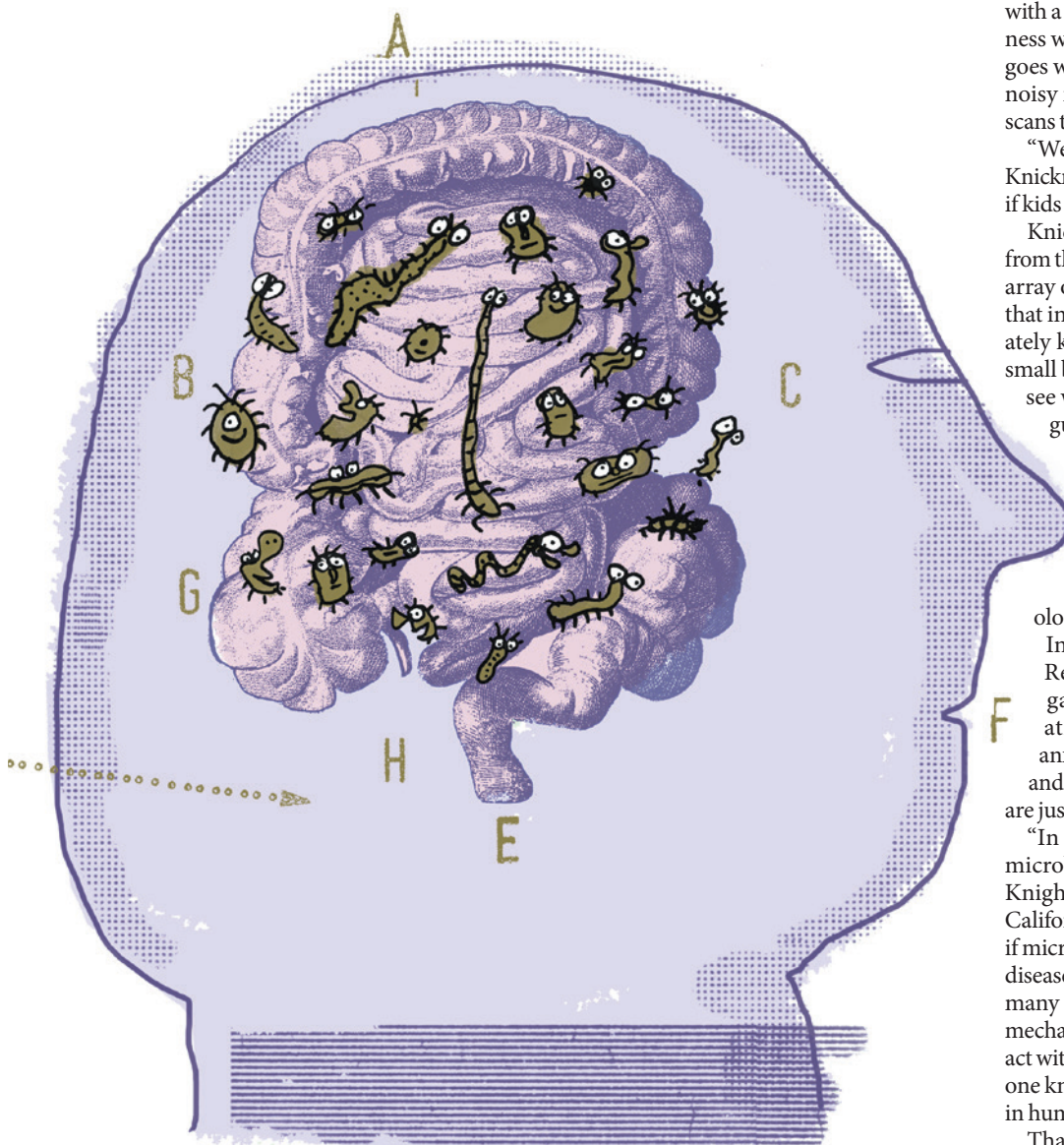


ILLUSTRATION BY SERGE BLOCH

the supplements industry from claiming that probiotics — bacteria that purportedly aid with digestive issues — can support emotional well-being. Pharmaceutical firms, hungry for new leads in treating neurological disorders, are beginning to invest in research related to gut microbes and the molecules that they produce.

Scientists and funders are looking for clarity. Over the past two years, the US National Institute of Mental Health (NIMH) in Bethesda, Maryland, has funded seven pilot studies with up to US\$1 million each to examine what it calls the ‘microbiome–gut–brain axis’ (Knickmeyer’s research is one of these studies). This year, the US Office of Naval Research in Arlington, Virginia, agreed to pump as much as US\$52 million over the next few years into work examining the gut’s role in cognitive function and stress responses. And the European Union has put €9 million (US\$10.1 million) towards a five-year project called MyNewGut, two main objectives of which target brain development and disorders.

The latest efforts aim to move beyond basic observations and correlations — but preliminary results hint at complex answers. Researchers are starting to uncover a vast, varied system in which gut microbes influence the brain through hormones, immune molecules and the specialized metabolites that they produce.

“There’s probably more speculation than hard data now,” Knickmeyer says. “So there’s a lot of open questions about the gold standard for methods you should be applying. It’s very exploratory.”

GUT REACTIONS

Microbes and the brain have rarely been thought to interact except in instances when pathogens penetrate the blood–brain barrier — the cellular fortress protecting the brain against infection and inflammation. When they do, they can have strong effects: the virus that causes rabies elicits aggression, agitation and even a fear of water. But for decades, the vast majority of the body’s natural array of microbes was largely uncharacterized, and the idea that it could influence neurobiology was hardly considered mainstream. That is slowly changing.

Studies on community outbreaks were one key to illuminating the possible connections. In 2000, a flood in the Canadian town of Walkerton contaminated the town’s drinking water with pathogens such as *Escherichia coli* and *Campylobacter jejuni*. About 2,300 people suffered from severe gastrointestinal infection, and many of them developed chronic irritable bowel syndrome (IBS) as a direct result.

During an eight-year study¹ of Walkerton residents, led by gastroenterologist Stephen Collins at McMaster University in Hamilton, Canada, researchers noticed that psychological issues such as depression and anxiety seemed to be a risk factor for persistent IBS. Premysl Bercik, another McMaster gastroenterologist, says that this interplay triggered

intriguing questions. Could psychiatric symptoms be driven by lingering inflammation, or perhaps by a microbiome thrown out of whack by infection?

The McMaster group began to look for answers in mice. In a 2011 study², the team transplanted gut microbiota between different strains of mice and showed that behavioural traits specific to one strain transmitted along with the microbiota. Bercik says, for example, that “relatively shy” mice would exhibit more exploratory behaviour when carrying the microbiota of more-adventurous mice. “I think it is surprising. The microbiota is really driving the behavioural phenotype of host. There’s a marked difference,” Bercik says. Unpublished research suggests that taking faecal bacteria from humans with both IBS and anxiety and transplanting it into mice induces anxiety-like

“THERE’S PROBABLY MORE SPECULATION THAN HARD DATA NOW.”

behaviour, whereas transplanting bacteria from healthy control humans does not.

Such results can be met with scepticism. As the field has developed, Knight says, microbiologists have had to learn from behavioural scientists that how animals are handled and caged can affect things such as social hierarchy, stress and even the microbiome.

And these experiments and others like them start with a fairly unnatural model: germ-free — or ‘gnotobiotic’ — mice. These animals are delivered by Caesarean section to prevent them from picking up microbes that reside in their mothers’ birth canals. They are then raised inside sterile isolators, on autoclaved food and filtered air. The animals are thus detached from many of the communal microbes that their species has evolved with for aeons.

In 2011, immunologist Sven Pettersson and neuroscientist Rochellys Diaz Heijtz, both at the Karolinska Institute in Stockholm, showed that in lab tests, germ-free mice demonstrated less-anxious behaviour than mice colonized with natural indigenous microbes³. (Less anxiety is not always a good thing, evolutionarily speaking, for a small mammal with many predators.)

When the Karolinska team examined the animals’ brains, they found that one region in germ-free mice, the striatum, had higher turnover of key neurochemicals that are associated with anxious behaviour, including the neurotransmitter serotonin. The study also showed that introducing adult germ-free mice to conventional, non-sterile environments failed to normalize their behaviour, but the offspring of such ‘conventionalized’ mice showed

some return to normal behaviour, suggesting that there is a critical window during which microbes have their strongest effects.

By this time, many researchers were intrigued by the mounting evidence, but results stemmed mostly from fields other than neuroscience. “The groups working on this are primarily gut folks, with a few psychology-focused people collaborating,” says Melanie Gareau, a physiologist at the University of California, Davis. “So the findings tended to describe peripheral and behavioural changes rather than changes to the central nervous system.”

But Pettersson and Diaz Heijtz’s research galvanized the field, suggesting that researchers could get past observational phenomenology and into the mechanisms affecting the brain. Nancy Desmond, a programme officer involved in grant review at the NIMH, says that the paper sparked interest at the funding agency soon after its publication and, in 2013, the NIMH formed a study section devoted to neuroscience research that aims to unravel functional mechanisms and develop drugs or non-invasive treatments for psychological disorders.

Judith Eisen, a neuroscientist at the University of Oregon in Eugene, earned a grant to study germ-free zebrafish, whose transparent embryos allow researchers to easily visualize developing brains. “Of course, ‘germ-free’ is a completely unnatural situation,” Eisen says. “But it provides the opportunity to learn which microbial functions are important for development of any specific organ or cell type.”

CHEMICAL EXPLORATION

Meanwhile, researchers were starting to uncover ways that bacteria in the gut might be able to get signals through to the brain. Pettersson and others revealed that in adult mice, microbial metabolites influence the basic physiology of the blood–brain barrier⁴. Gut microbes break down complex carbohydrates into short-chain fatty acids with an array of effects: the fatty acid butyrate, for example, fortifies the blood–brain barrier by tightening connections between cells (see “The gut–brain axis”).

Recent studies also demonstrate that gut microbes directly alter neurotransmitter levels, which may enable them to communicate with neurons. For example, Elaine Hsiao, a biologist now at the University of California, Los Angeles, published research⁵ this year examining how certain metabolites from gut microbes promote serotonin production in the cells lining the colon — an intriguing finding given that some antidepressant drugs work by promoting serotonin at the junctions between neurons. These cells account for 60% of peripheral serotonin in mice and more than 90% in humans.

Like the Karolinska group, Hsiao found that germ-free mice have significantly less serotonin floating around in their blood, and she also showed that levels could be restored by introducing to their guts spore-forming bacteria (dominated by *Clostridium*, which break down

short-chain fatty acids). Conversely, mice with natural microbiota, when given antibiotics, had reduced serotonin production. “At least with those manipulations, it’s quite clear there’s a cause–effect relationship,” Hsiao says.

But it remains unclear whether these altered serotonin levels in the gut trigger a cascade of molecular events, which in turn affect brain activity — and whether similar events take place in humans, too. “It will be important to replicate previous findings, and translate these findings into human conditions to really make it to the textbooks,” Hsiao says.

For John Cryan, a neuroscientist at University College Cork in Ireland, there is little question that they will. His lab has demonstrated⁶ that germ-free mice grow more neurons in a specific brain region as adults than do conventional mice. He has been promoting the gut–brain axis to neuroscientists, psychiatric-drug researchers and the public. “If you look at the hard neuroscience that has emerged in the last year alone, all the fundamental processes that neuroscientists spend their lives working on are now all shown to be regulated by microbes,” he says, pointing to research on the regulation of the blood–brain barrier, neurogenesis in mice and the activation of microglia, the immune-like cells that reside in the brain and spinal cord.

At the 2015 Society for Neuroscience meeting in Chicago, Illinois, this month, Cryan and his colleagues plan to present research showing that myelination — the formation of fatty sheathing that insulates nerve fibres — can also be influenced by gut microbes, at least in a specific part of the brain. Unrelated work⁷ has shown that germ-free mice are protected from an experimentally induced condition similar to multiple sclerosis, which is characterized by demyelination of nerve fibres. At least one company, Symbiotix Biotherapies in Boston, Massachusetts, is already investigating whether a metabolite produced by certain types of gut bacterium might one day be used to stem the damage in humans with multiple sclerosis.

A MOVE TO THERAPY

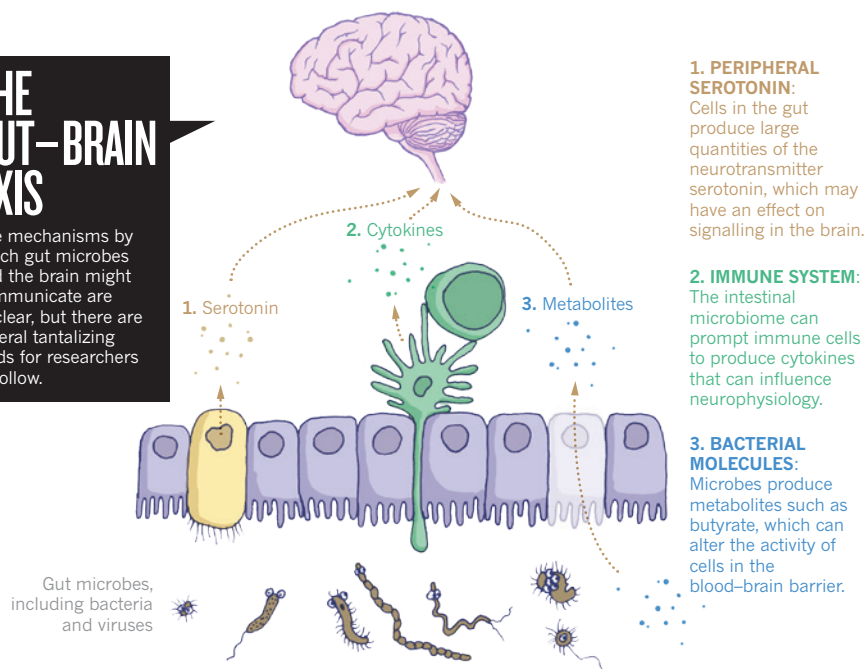
Tracy Bale, a neuroscientist at the University of Pennsylvania in Philadelphia, suspects that simple human interventions may already be warranted. Bale heard about Cryan’s work on the radio programme *Radiolab* three years ago. At the time, she was researching the placenta, but wondered how microbes might fit into a model of how maternal stress affects offspring.

In research published this year⁸, Bale subjected pregnant mice to stressful stimuli. She found that it noticeably reduced the levels of *Lactobacilli* present in the mice’s vaginas, which are the main source of the microbes that colonize the guts of offspring. These microbial shifts carried over to pups born vaginally, and Bale detected signs that microbiota might affect neurodevelopment, especially in males.

In work that her group plans to present at the Society for Neuroscience meeting, Bale has

THE GUT–BRAIN AXIS

The mechanisms by which gut microbes and the brain might communicate are unclear, but there are several tantalizing leads for researchers to follow.



shown that by feeding vaginal microbiota from stressed mice to Caesarean-born infant mice, they can recapitulate the neurodevelopmental effects of having a stressed mother. Bale and her colleagues are now wrapping up research investigating whether they can treat mice from stressed mums with the vaginal microbiota of non-stressed mice.

The work, Bale says, has “immediate translational effects”. She points to a project headed by Maria Dominguez-Bello, a microbiologist at the New York University School of Medicine, in which babies born by means of Caesarean section are swabbed on the mouth and skin with gauze taken from their mothers’ vaginas. Her team wants to see whether these offspring end up with microbiota similar to babies born vaginally. “It’s not standard of care,” Bale says, “but I will bet you, one day, it will be.”

Many are still sceptical about the link between microbes and behaviour and whether it will prove important in human health — but scientists seem more inclined to entertain the idea now than they have in the past. In 2007, for example, Francis Collins, now director of the US National Institutes of Health, suggested that the Human Microbiome Project, a large-scale study of the microbes that colonize humans, might help to unravel mental-health disorders. “It did surprise a few people who assumed we were talking about things that are more intestinal than cerebral,” Collins says. “It was a little bit of leap, but it’s been tentatively backed up.”

Funding agencies are supporting the emerging field, which spans immunology, microbiology and neuroscience, among other disciplines. The NIMH has offered seed funding for work on model systems and in humans to probe whether the area is worth more-substantial investment, a move that has already brought more researchers into the fold. The MyNewGut project in Europe has an even more optimistic view of the value of such research, specifically

seeking concrete dietary recommendations that might alleviate brain-related disorders.

Today, Knickmeyer’s project on infants represents what she calls “a messy take-all-comers kind of sample”. Among the brain regions that Knickmeyer is scanning, the amygdala and prefrontal cortex hold her highest interest; both have been affected by microbiota manipulations in rodent models. But putting these data together with the dozens of other infant measures that she is taking will be a challenge. “The big question is how you deal with all the confounding factors.” The children’s diets, home lives and other environmental exposures can all affect their microbiota and their neurological development, and must be teased apart.

Knickmeyer speculates that tinkering with microbes in the human gut to treat mental-health disorders could fail for other reasons. Take, for instance, how microbes might interact with the human genome. Even if scientists were to find the therapeutic version of a “gold Cadillac of microbiota”, she points out, “maybe your body rejects that and goes back to baseline because your own genes promote certain types of bacteria.” There is much more to unravel, she says. “I’m always surprised. It’s very open. It’s a little like a Wild West out there.” ■

Peter Andrey Smith is a reporter based in New York City.

1. Marshall, J. K. *et al.* *Gut* **59**, 605–611 (2010).
2. Bercik, P. *et al.* *Gastroenterology* **141**, 599–609 (2011).
3. Diaz Heijtz, R. *et al.* *Proc. Natl Acad. Sci. USA* **108**, 3047–3052 (2011).
4. Braniste, V. *et al.* *Sci. Transl. Med.* **6**, 263ra158 (2014).
5. Yano, J. M. *et al.* *Cell* **161**, 264–276 (2015).
6. Ogbonaya, E. S. *et al.* *Biol. Psychiatry* **78**, e7–e9 (2015).
7. Lee, Y.-K., Menezes, J. S., Umesaki, Y. & Mazmanian, S. K. *Proc. Natl Acad. Sci. USA* **108**, 4615–4622 (2010).
8. Jašarević, E., Howerton, C. L., Howard, C. D. & Bale, T. L. *Endocrinology* **156**, 3265–3276 (2015).



COMMENT

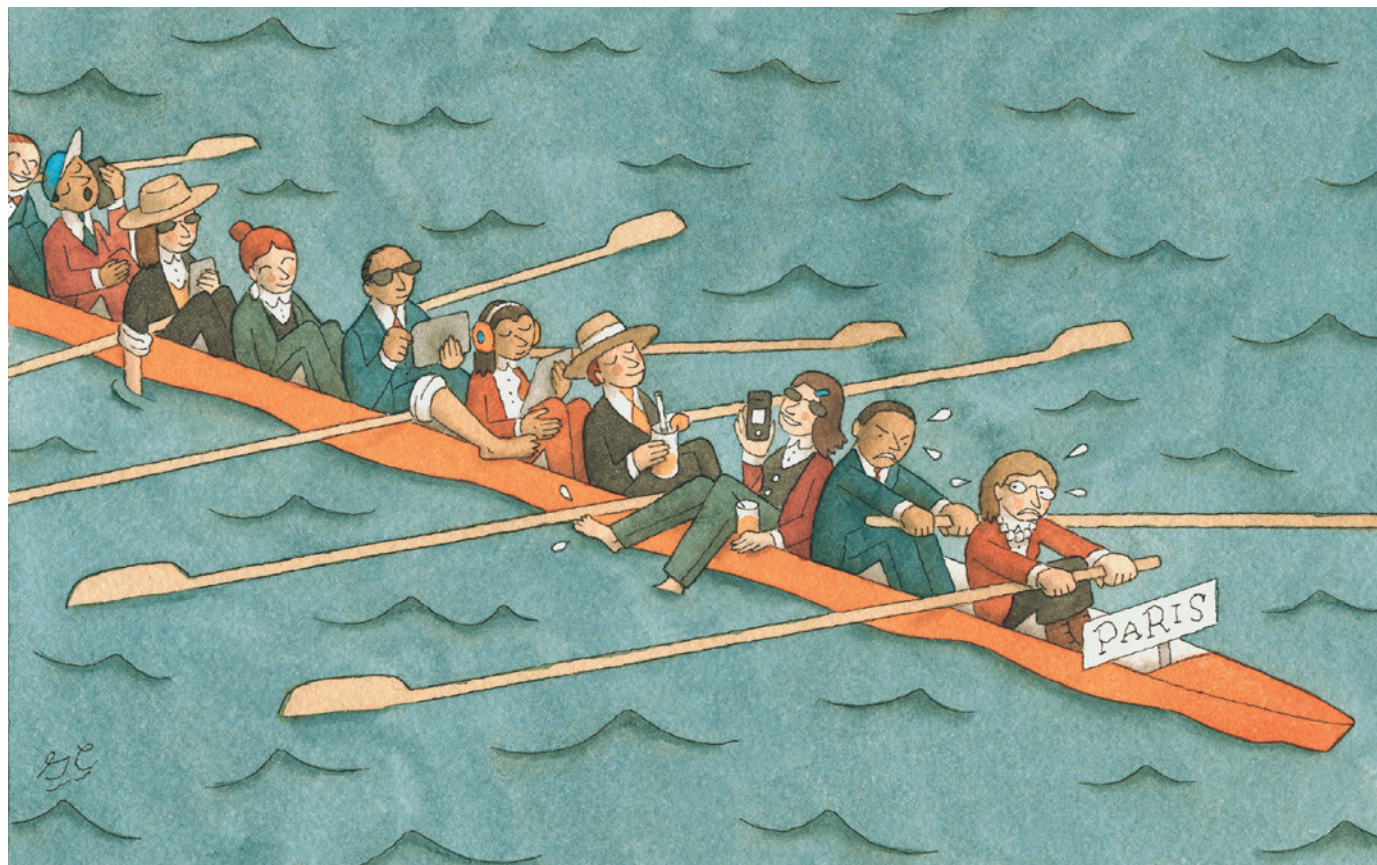
POLICY Eight ways to get the best out of experts **p.317**

AI Three books agree that human psychology is a barrier to robot superiority **p.320**

MUSIC In conversation with the creator of a chamber opera exploring time **p.322**

OBITUARY William Paul, discoverer of interleukin-4, remembered **p.324**

ILLUSTRATION BY GREG CLARKE



Price carbon — I will if you will

To forge a strong climate accord in Paris, nations must agree on a common goal in everyone's self-interest, say **David J. C. MacKay** and colleagues.

Negotiations at the United Nations climate summit in Paris this December will adopt a 'pledge and review' approach to cutting global carbon emissions. Countries will promise to reduce their emissions by amounts that will be revised later. The narrative is that this will "enable an upward spiral of ambition over time"¹. History and the science of cooperation predict that quite the opposite will happen.

Climate change is a serious challenge because the atmosphere gives a free ride to countries that emit. If some nations sit back and rely on others' efforts, the incentives for anyone to act are weakened. Review of the first phase of the Kyoto Protocol at the 2012 UN climate meeting in Doha, for instance, resulted in Japan, Russia, Canada and New Zealand leaving the agreement, frustrating those who kept their promises.

Success requires a common commitment, not a patchwork of individual ones. Negotiations need to be designed to realign self-interests and promote cooperation. A common commitment can assure participants that others will match their efforts and not free-ride. A strategy of "I will if you will" stabilizes higher levels of cooperation. It is the most robust pattern of cooperation seen in laboratory and field studies of ►

► situations open to free-riding².

A global carbon price — so far excluded from consideration in international negotiations — would be the ideal basis for a common commitment in our view. A price is easy to agree and handle, relatively fair, less vulnerable to gaming than global cap-and-trade systems, and consistent with climate policies already in place, such as fossil-fuel taxes and emissions cap-and-trade.

Only a common commitment can lead to a strong treaty. Forty years of empirical and theoretical literature on cooperation confirms that individual commitments do not deliver strong collective action. Cooperators find that defectors take advantage of them. Ambition declines when others are revealed to be free-riding³. Dishes often stack up in the sinks of shared apartments. But in the Alps, villagers have successfully managed shared land for hundreds of years, with a common commitment governing grasslands⁴.

COMMON COMMITMENT

Imagine that you and nine other self-interested players (representing countries) take part in a game. Each player has \$10, some or all of which they may simultaneously pledge to a common pot. A referee makes sure that they honour their pledges. Every dollar (for carbon dioxide abatement) placed in the pot will be doubled (by climate benefits) and distributed evenly to all players. So putting a dollar in the pot will return 20 cents to each player.

Consider two variants of the game. First, in the 'individual commitment' version, pledges are made independently. This is the classic public-goods game, in which the rational selfish strategy is to contribute nothing, because this makes a player better off no matter what the others do. The result is the famous tragedy of the commons. Cooperation does not occur, even though everyone would gain from it.

Second, in the 'common commitment' version, players condition their contributions on others' pledges: a referee ensures that all contribute the amount of the lowest pledge. After enforcing this common commitment, the money is doubled and distributed evenly, exactly as before.

This changes everything. Pledging \$0 will mean simply keeping your \$10, whereas pledging \$10 could result in ending up with anything between \$10 and \$20, depending on what others pledge. So, because you cannot lose and could gain by pledging \$10, that is what you would do, even if you are completely selfish. Since all parties would pledge \$10, the group's \$100 is doubled and all end up with the maximum amount of \$20.

Selfish behaviour has been changed from 'contribute nothing' to 'contribute everything', because the common commitment

protects against free-riding.

In 1997, the Kyoto negotiators initially did try to agree a common commitment, expressed as a formula for national emissions caps, but failed. In the end, each nation was simply asked to submit their final numbers for insertion into the draft annex⁵. The result was a patchwork of weak and unstable commitments. Similarly, in response to the 2009 Copenhagen Accord, China pledged emissions equal to those considered 'business as usual' before the accord; and India pledged even less.

Enforcement is widely thought to be the missing ingredient in the Kyoto Protocol and crucial for the success of a Paris agreement. This is only half right — both enforcement and a common commitment are required. For example, if drivers chose their own speed limits, there would be no use enforcing them, because everyone would drive at their desired speed. Instead, because it limits others as well, people agree to a common speed limit that is lower than almost everyone's individual limit. In other words, with individual commitments, there is nothing meaningful to enforce, whereas enforcement strengthens a common commitment.

What could all countries commit to? National limits on the quantity of emissions will not work. Kyoto negotiators suggested at least ten formulae to determine the reductions that each nation should make, but could not agree. When attention turned to reducing

"Harness self-interest by aligning it with the common good."

emissions by some percentage relative to 1990 levels, individual commitments ranged from an 8% decrease to a 10% increase. The United States and developing countries

made no commitments at all.

Percentage pledges failed because countries differ; for instance, some economies declined after 1990 and some grew. Developing countries fear caps that curb their growth. Instead they see it as fair to allocate emission permits on an equal per capita basis. Because permit sales would result in huge wealth transfers to poor countries, rich countries find such proposals unacceptable⁶.

There is no longer any serious discussion of a common commitment to reduce the quantity of carbon emissions.

GLOBAL CARBON PRICE

We, and others, propose an alternative: a global carbon-price commitment⁷. Each country would commit to place charges on carbon emissions from fossil-fuel use (by taxes or cap-and-trade schemes, for example) sufficient to match an agreed global price, which could be set by voting — by a super-majority rule that would produce a coalition of the willing.

A uniform carbon price is widely accepted as the most cost-effective way to curb emissions. Carbon pricing is flexible, allowing fossil taxes, cap-and-trade, hybrid schemes and other national policies to be used (unlike a global carbon tax). All that is required of a country is that its average carbon price — cost per unit of greenhouse gas emitted — be at least as high as the agreed global carbon price.

Unlike global cap-and-trade, carbon pricing allows countries to keep all carbon revenues, eliminating the risk of needing to buy expensive credits from a rival country. Taxes need not rise if a nation performs a green tax shift — reducing taxes on good things such as employment by charging for pollution. Shifting taxes from good things to bad things could mean there is no net social cost to pricing carbon, even before counting climate benefits⁸.

A global price does not automatically result in acceptable burden sharing. A 'Green Climate Fund' will be needed to transfer funds from rich to poor countries. To minimize disputes, the objective of climate-fund transfers should be to maximize the global price of carbon. This can be implemented in a way that encourages rich countries to be generous and poor countries to vote for a higher global carbon price⁹, for example, by making all climate-fund payments proportional to the agreed carbon price.

After decades of failure, a fresh approach is needed — one that is guided by the science of cooperation. A common price commitment would harness self-interest by aligning it with the common good. Nothing could be more fundamental. ■

David J. C. MacKay is professor of engineering at the University of Cambridge, UK. **Peter Cramton** is professor of economics at the University of Maryland, College Park, Maryland, USA. **Axel Ockenfels** is professor of economics at the University of Cologne, Germany. **Steven Stoft** is an economist in Berkeley, California, USA.
e-mail: steven@stoft.com

1. Ad Hoc Working Group on the Durban Platform. *Parties' Views and Proposals on the Elements for a Draft Negotiating Text ADP.2014.6* (UNFCCC, 2014); available at <http://go.nature.com/x1fjcd>
2. Kraft-Todd, G., Yoeli, E., Bhanot, S. & Rand, D. *et al. Curr. Op. Behav. Sci.* **3**, 96–101 (2015).
3. Ledyard, J. in *The Handbook of Experimental Economics* (eds Kagel, H. & Roth A. E.) 111–194 (Princeton Univ. Press, 1995).
4. Ostrom, E. *Governing the Commons: The Evolution of Institutions for Collective Action* (Cambridge Univ. Press, 1990).
5. Depledge, J. *The Origins of the Kyoto Protocol* (UNFCCC, 2000).
6. Stiglitz, J. in *Making Globalization Work* Ch. 6 (Norton, 2006).
7. Cramton, P., Ockenfels, A. & Stoft, S. *Econ. Energy Environ. Policy* **4**, 1–4 (2015).
8. Bovenberg, A. L. *Int. Tax Public Finance* **6**, 421–443 (1999).
9. Cramton, P. & Stoft, S. *Econ. Energy Environ. Policy* <http://dx.doi.org/10.5547/2160-5890.1.2.9> (2012).



Use experts wisely

Policymakers are ignoring evidence on how advisers make judgements and predictions, warn
William J. Sutherland and Mark A. Burgman.

Many governments aspire to evidence-based policy and practice. The predominant, conventional approaches to using experts are either to seek the advice of one highly regarded individual, or to convene a panel with diverse expertise across relevant areas. For example, quarantine services worldwide routinely rely on expert judgement to estimate the probability of entry, establishment and spread of pests and diseases.

The accuracy and reliability of expert opinions is, however, compromised by a long list

of cognitive frailties¹. Estimates are influenced by experts' values, mood, whether they stand to gain or lose from a decision², and by the context in which their opinions are sought. Experts are typically unaware of these subjective influences. They are often highly credible, yet they vastly overestimate their own objectivity and the reliability of their peers³.

Happily, a large and growing body of literature describes methods for engaging with experts that enhance the accuracy and calibration of their judgements^{4,5}. Unhappily, these methods are rarely used to support

public-policy decisions. All the methods strive to alleviate the effects of psychological and motivational bias; all structure the acquisition of estimates and associated uncertainties; and all recommend combining independent opinions. None relies on the opinion of the best-regarded expert or uses unstructured group consensus.

The cost of ignoring these techniques — of using experts inexpertly — is less-accurate information, and thus more frequent and more serious policy failures.

KNOWN AND UNKNOWN

For an important subset of questions, expert technical judgements about facts plays a part in policy and decision-making. (We appreciate that political context may determine what comprises relevant, convincing evidence, and that that evidence rarely leads directly to policy and action because decision-makers must balance a range of political, social, economic, practical and scientific issues.)

Policymakers use expert evidence as though it were data. So they should treat expert estimates with the same critical rigour that must be applied to data. Experts must be tested, their bias minimized, their accuracy improved, and their estimates validated with independent evidence (see 'Eight ways to improve expert advice'). That is, experts should be held accountable for their opinions.

For example, experts who are confident and routinely close to the correct answer provide more information than do experts who regularly deviate from the correct answer or are under-confident. Highly regarded experts are routinely shown to be no better than novices at making judgements. Opinions from more-informative experts can be weighted more heavily, whereas the opinions of some experts may be discarded altogether⁶. These strategies will illuminate where advice is robust, and where it is contradictory, self-serving or misguided. This will generate evidence for policy decisions that is more relevant and reliable. Roger Cooke, a risk-analysis researcher at the Delft University of Technology in the Netherlands and his colleagues have used this approach effectively to better predict the implications of policy for transport and nuclear-power safety⁴.

Experts themselves must make explicit the sensitivity of their decisions to scientific uncertainty, assumptions and caveats. When invited to advise, they should demand that state-of-the-art techniques are used to harvest and process what they offer. If not, all involved risk wasting substantial time, resources and opportunities.

Importantly, all parties must be clear about what they expect of each other: estimates of facts, predictions of event outcomes or advice on the best course of action. Properly managed, experts can help with the

first two. Providing advice assumes that the expert shares the same values and objectives as the decision-makers.

Several processes have been shown to improve experts' performances on estimates of facts and predictions of event outcomes. In using specialists to weigh up the best course of action, the researchers themselves, and the policymakers using them, should identify all possible changes, options and threats, a process known as horizon scanning. Policymakers must list all possible known solutions using a wide group of experts and reference to the literature, to reduce the risk that valuable alternatives are overlooked (known as solution scanning⁷).

Deliberations should be underpinned by a systematic collection of evidence, an assessment of its relevance, and an identification of the knowledge gaps that might change the decision. The information can be collated so that it is ready for use — rather than in response to a policy need — as is being done for biodiversity⁸ (see www.conservationevidence.com).

RULES OF ENGAGEMENT

A few more rules of engagement, routinely applied, will enhance the quality and reliability of expert judgements.

Ensure that questions are fully specified and unambiguous, so that language-based uncertainties do not cloud judgements. For example, a seemingly straightforward question such as 'How many diseased animals are there in the region?' could be interpreted differently by different people. The question

does not specify whether to include only those animals that are known to be infectious, or also those that have died, have recovered, are diseased but yet to be identified as such, and so on.

Structured question formats counter tendencies towards over-confidence for individual estimates. For example, Andrew Speirs-Bridge at La Trobe University has shown⁹ that questions that elicit four responses — upper and lower bounds, a best guess and a degree of confidence in the interval — generate estimates that are relatively well calibrated. Consider a range of scenarios and alternative theories. Ensure that several experts answer each question.

Unstructured group interactions are subject to 'groupthink': the group gravitates towards an initial or even an arbitrary estimate; dominant individuals drive the outcome; or individuals are ascribed greater credibility than they deserve because of their appearance, manner or professional background. Structured, facilitated interactions counter factors such as these, which distort estimates³.

Review assumptions, reconcile misunderstandings and introduce new information. Ensure that decision-makers do not rely on experts to choose between options but rather use an appropriate decision tool. One such is structured decision-making, in

which experts populate decision tables with estimates of the expected outcomes for each criterion under each policy option, but do not decide the best option.

For example, an analysis⁶ of volcano-eruption risks by Willy Aspinall, an Earth scientist at the University of Bristol, UK, used structured interactions. These substantially improved the quality of estimates, because he ensured that well-specified questions were answered by several experts in such a way that he avoided or mitigated the psychological tripwires that compromise many group interactions.

Similarly, a study led by conservation ecologist Marissa McBride¹⁰ at the University of Melbourne in Australia engaged with groups of experts remotely, using structured questions and group interactions to assess the conservation status of threatened Australian birds. They used telephone conferences to outline the context and purpose of the interactions, which was to reassess the International Union for Conservation of Nature's Red List assessments for a suite of threatened species. They then used e-mail to elicit initial judgements and to clarify questions, introduce further data and explanations. Finally, they circulated a spreadsheet and compiled a second round of private, anonymous judgements.

In many cases, incorporating the formal stages described here will improve decision-making. The benefits are substantial improvements in the reliability of judgements, relatively free of personal biases and values. The costs in time and resources are modest. ■

"Structured question formats counter tendencies towards over-confidence"

LEVERAGE THE LITERATURE

Eight ways to improve expert advice

Use groups. Their estimates consistently outperform those of individuals.

Choose members carefully. Expertise declines dramatically outside an individual's specialization or experience.

Don't be starstruck. Age, number of publications, technical qualifications, years of experience, memberships of learned societies and apparent impartiality do not explain an expert's ability to estimate unknown quantities or predict events. This finding applies in studies from nuclear-safety systems and geopolitics to ecology.

Avoid homogeneity. Diverse groups tend to generate more-accurate judgements.

Don't be bullied. People who are less self-assured and assertive, and who integrate

information from diverse sources tend to make better judgements.

Weight opinions. Calibrate experts' performance with test questions. This improves risk estimates in many domains, including earthquakes and nuclear-safety systems.

Train experts. Training can improve experts' abilities to estimate probabilities of events, quantities or model parameters.

Give feedback. Chess players, weather forecasters, sports people, gamblers, intensive-care physicians and physicists solving textbook problems generally make accurate judgements, probably as a result of rapid feedback from mistakes that are visible and personal. Experts deserve the same — give them immediate and unambiguous feedback.

William J. Sutherland is professor of conservation biology in the Department of Zoology, University of Cambridge, UK.

Mark Burgman is professor of botany and managing director of the Centre of Excellence for Biosecurity Risk Analysis, School of BioSciences, University of Melbourne, Parkville, Australia.
e-mail: wjs32@cam.ac.uk

1. Tversky, A. & Kahneman, D. in *Judgement Under Uncertainty: Heuristics and Biases* (eds Kahneman, D., Slovic, P. & Tversky, A.) 23–30 (Cambridge Univ. Press, 1982).
2. Englich, B. & Soder, K. *Judgm. Decis. Mak.* **4**, 41–50 (2009).
3. Burgman, M. A. et al. *PLoS ONE* **6**, e22998 (2011).
4. Cooke, R. M. *Experts in Uncertainty: Opinion and Subjective Probability in Science* (Oxford Univ. Press, 1991).
5. O'Hagan, A. et al. *Uncertain Judgements: Eliciting Experts' Probabilities* (Wiley, 2006).
6. Aspinall, A. *Nature* **463**, 294–295 (2010).
7. Sutherland, W. J. et al. *Ecol. Soc.* **19**, 3 (2014).
8. Sutherland, W. J. et al. *What Works in Conservation* (OpenBooks, 2015).
9. Speirs-Bridge, A. et al. *Risk Analysis* **30**, 512–523 (2010).
10. McBride, M. F. et al. *Methods Ecol. Evol.* **3**, 906–920 (2012).

For a list of further reading on this topic, see go.nature.com/ibrrqd.



Fear of robots overlooks their limitations, and potential.

ROBOTICS

Countering singularity sensationalism

Ken Goldberg reviews three books that probe the nexus of people and robots.

In the late nineteenth century, the United States was awash with the racist term 'yellow peril'. Fears spread that Chinese immigrants working in the country's mines and building its railways would seize more jobs from the citizenry. Today, there is a similar collective fear, this time about a 'singularity' in which artificial intelligence (AI) and robots surpass human abilities. In May 2014, for instance, physicists Stephen Hawking, Frank Wilczek and Max Tegmark, with computer scientist Stuart Russell, warned in the newspaper *The Independent*: "Success in creating AI would be the biggest event in human history. Unfortunately, it might also be the last".

Surprising advances are being achieved, for example in 'deep learning' — a method for approximating complex functions using thousands of numerical parameters. And robots are evolving, with advances in 3D sensing and mapping. But progress is not nearly as steady as some claim. Three books explore the topic from different perspectives. All suggest that robot superiority faces a formidable obstacle: human psychology.

In *Machines of Loving Grace*, reporter John Markoff of *The New York Times* highlights

Machines of Loving Grace: The Quest for Common Ground Between Humans and Robots

JOHN MARKOFF
Ecco: 2015.

Rise of the Robots: Technology and the Threat of a Jobless Future

MARTIN FORD
Basic: 2015.

Our Robots, Ourselves: Robotics and the Myths of Autonomy

DAVID A. MINDELL
Viking: 2015.

the compelling contrast between AI and intelligence amplification (IA). He chronicles the fascinating and often antagonistic evolution of these fields since 1956, when both terms were coined.

Markoff has been interacting with leading researchers for the past 20 years. As he shows, despite early optimism, creating AI has turned out to be extremely difficult. Robotics remains challenged by Moravec's paradox: tasks that are tough for humans, such as precision spot welding, are easy for robots, whereas tasks that are easy for humans, like reliably clearing a dinner table,

remain extremely hard for robots. This is mostly attributable to the inherent complexity of friction, collisions and contact mechanics. It is much easier to calculate the precise trajectory of a comet than to predict that of a coffee mug pushed across a tabletop.

The term IA was preferred by Douglas Engelbart, who worked on the potential of computers to amplify human abilities. This evolved into the field of human-computer interaction, which brought us the mouse and graphical interfaces. Markoff recounts the stories of pioneers such as computer scientist Terry Winograd and augmented-reality expert Gary Bradski. Both recognized the limitations of AI, and became advocates of IA. Markoff's book makes a strong case that the success of AI will depend on advances in IA.

Rise of the Robots by software entrepreneur Martin Ford proclaims that AI and robots are about to eliminate most jobs, blue- and white-collar. A close reading reveals the evidence as extremely sketchy. Ford has swallowed the rhetoric of futurist Ray Kurzweil, and repeatedly asserts that we are on the brink of vastly accelerating advances based on Moore's law, which posits that computing power increases exponentially with time. Yet some computer scientists rue this exponential fallacy, arguing that the success of integrated circuits has raised expectations of progress far beyond what historians of technology recognize as an inevitable flattening of the growth curve.

Nor do historical trends support the Luddite fallacy, which assumes that there is a fixed lump of work and that technology inexorably creates unemployment. Such reasoning fails to consider compensation effects that create new jobs, or myriad relevant factors such as globalization and the democratization of the workforce. Ford describes software systems that attempt to do the work of attorneys, project managers, journalists, computer programmers, inventors and musicians. But his evidence that these will soon be perfected and force massive lay-offs consists mostly of popular magazine articles and, in one case, a conversation with the marketing director of a start-up.

In *Our Robots, Ourselves*, telerobotist David Mindell points out that autonomous systems are not new. Since the 1970s, they have been in daily use at very low and very high altitudes, for deep-sea and space exploration and in almost all aircraft. Drawing on extensive experience, Mindell explains that although such systems have evolved, many experts continue to mistrust them. For example, there is a persistent school of thought that oceanographers must experience the murky depths directly to understand the wonders lurking there. Yet robotic submarines, tele-operated through fibre-optic cables, are much more agile, are able to explore for longer and do not require expensive certification for every modification.

Mindell recounts a telling scene. In 1977, oceanographer Robert Ballard of the Woods Hole Oceanographic Institution in Massachusetts and microbiologist Holger Jannasch descended in a submersible. They came across a sea-floor vent field. As Ballard looked through the glass portal at the “crab gradient” near the fissure’s source, he realized that Jannasch had his back to it. When asked why, Jannasch replied that the view was “better here”, indicating the television image relayed from the craft’s camera. In an ‘aha’ moment, Ballard realized that the view will always be better at the surface, where high-resolution video from below can be viewed in comfort. Yet many continue to argue otherwise.

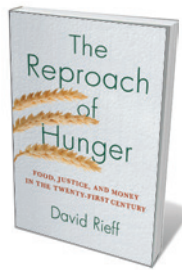
Mindell looks at the compelling historical, cultural, political, psychological, philosophical and public-relations justifications for keeping people in control. A recent example is the July 2015 international petition to ban autonomous weapons, signed by almost 3,000 researchers. Consider, for instance, Google’s announcement that it had found the most unreliable component in its autonomous cars: human drivers trying to take over. Its response was to remove the steering wheel. Mindell explains why this is a mistake. In what he calls the “myth of full autonomy”, he points out that a machine may operate on its own for intervals, but that no machine works entirely independently: human intentions, assumptions and parameters are built into all machines. Mindell’s experience leads him to a similar conclusion to Markoff’s: the essential (and most difficult) challenge is designing the interfaces that keep humans in the loop.

Technological progress is not deterministic. There have been several cycles of irrational enthusiasm followed by disappointment and the evaporation of research and development funding, known as AI winters. The latest round of expectations feels similarly exaggerated.

Alarmist writing may hasten the next slump and distract attention from a more realistic and important development, which we might call multiplicity. Multiplicity characterizes an emerging category of systems in which diverse groups of humans work together with diverse groups of machines to solve difficult problems, consistent with the points made by Markoff and Mindell. Multiplicity thrives at the intersection of AI and IA, combining the wisdom of crowds with the power of cloud computing. As designer and computer scientist John Maeda has put it, it is not us versus the machines; it is us and the machines. There is much to gain by joining forces. ■

Ken Goldberg is professor of engineering and director of the People and Robots Initiative at the University of California, Berkeley.
e-mail: goldberg@berkeley.edu

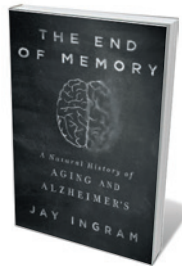
Books in brief



The Reproach of Hunger: Food, Justice and Money in the 21st Century

David Rieff VERSO (2015)

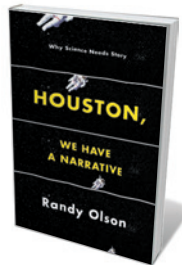
As refugee crises fill the news, David Rieff reminds that hunger is a war not won. Rieff, a veteran thinker on development issues, spent six years researching the nexus of population, food commodification and persistent poverty for this critical analysis. Scathing about the alarmist or over-optimistic pronouncements of development officials, agribusiness multinationals and philanthropic nabobs, he notes that any issue involving billions of humans cannot be neatly engineered. Thoughtful, trenchant and bracingly sceptical.



The End of Memory: A Natural History of Aging and Alzheimer's

Jay Ingram THOMAS DUNNE (2015)

Alzheimer's disease affects 5.3 million US citizens, and has so far eluded cure. In this deft overview, science writer Jay Ingram unravels the complexities of the science past and present. He examines the legacy of neurology pioneers such as Aloisius Alzheimer and Frederic Lewy; the biology of ageing and shifts in episodic and autobiographical memory; and the protein plaques and neurofibrillary tangles associated with the disease. And there is more, from the ongoing Nun Study of Aging and Alzheimer's Disease (begun by neurologist David Snowden in 1986) to the idea of Alzheimer's as “type 3 diabetes”.



Houston, We Have a Narrative: Why Science Needs Story

Randy Olson UNIVERSITY OF CHICAGO PRESS (2015)

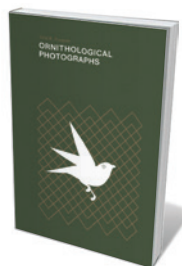
Whether synthetic biology or exoplanet hunting, science told well can carry a thriller-like punch. Marine biologist turned filmmaker Randy Olson argues that narrative skill is central not just to science communication but also to research reportage, preventing false positives, yawn-worthy delivery and more. Olson prescribes the Hollywood formula “and, but, therefore” as the backbone of story, introducing momentum, conflict and resolution. He has packed his solid primer with analyses of how it is done, from James Watson's 1968 *The Double Helix* (Atheneum) to exemplary scientific abstracts.



Memory and Movies: What Films Can Teach Us about Memory

John Seamon MIT PRESS (2015)

Cinema has long exploited the dramatic potential of memory. Here, John Seamon exploits film's potential for elucidating neuroscience. Inspired by Christopher Nolan's 2000 *Memento* (which hinges on anterograde amnesia), Seamon trains a cinematic lens on aspects of memory from facial recognition to dementia. Philip Kaufman's 1978 *Invasion of the Body Snatchers*, for instance, mirrors Capgras' delusion, in which people believe that their acquaintances are doppelgängers; while Robert Redford's *Ordinary People* (1980) dissects post-traumatic stress disorder with exquisite precision.



Ornithological Photographs

Todd Forsgren DAYLIGHT (2015)

Photographer Todd Forsgren has spent years capturing images of birds, from hummingbirds to toucans, caught in mist nets — a tool widely used by ornithologists for ring-banding and data collection. Some may find the sight of immobilized birds in this collection disturbing. But Forsgren's book uniquely showcases the birds' individuality while testifying to the painstaking, ongoing work of field researchers striving to understand the ecology, population flux and more of wild birds. [Barbara Kiser](#)



Refuse the Hour explores the unshakable anxiety of disappearing forever.

Q&A Peter Galison

Time transformer

Next week, *Refuse the Hour*, a chamber opera about time, opens at the Brooklyn Academy of Music in New York City. The work is a collaboration between physics historian Peter Galison and South African multimedia artist William Kentridge. Galison talks about the nexus of technology and imperial conquest, the 'twin paradox' associated with Einstein's special theory of relativity and the metaphorical resonance between black holes and mortality.



How did this collaboration begin?

I met William Kentridge through a mutual friend who was interested in fostering collaboration between science and art. What struck me was the probing psychological insight that William brought to the history of empire, how he depicted colonialism and its technical culture, surrounded by the terror of the historical moment. We began to focus on time. Early on, William said, "I don't want to make an illustrated science lecture." I said, "That's great, because I don't want to be the science adviser to an art project." Instead, we kept returning to scientific arguments about the nature of time and using them as jumping-off points — trying to capture, in affect and aesthetics, a sense of what gave them force.

What will audiences see in *Refuse the Hour*?

We originally had in mind an elaborate combination of installation and performance, but we then split the project into a five-channel video installation called *The Refusal of Time*, and an 80-minute chamber opera, *Refuse the Hour*. The opera is about the transformations that we have witnessed in our cultural

and scientific ideas of time: Isaac Newton's absolute time, Albert Einstein's relative time and the threatened end to time and space as one approaches the centre of a black hole. When the opera makes its North American premiere, the audience will see projections of giant metronomes, African singers with megaphones, a giant accordion-like 'breathing machine' and William lecturing on procrastination, entropy and empire.

How were you inspired by historical timekeeping?

At the turn of the twentieth century, to measure longitude well enough to map the world's coastlines, colonial powers such as Britain needed to coordinate timekeeping across the globe. In a period when most countries lacked electric light, these powers were stringing telegraph cables under the ocean from Europe to Africa and South America. It was an extraordinary moment of scientific ambition and brute imperialism.

Was there a backlash against the standardization of time?

The opera shows both hope and resistance associated with advances in timekeeping. In the 1870s, authorities in Paris and Vienna

Refuse the Hour
22–25 October 2015.
Brooklyn Academy of
Music, New York City.

tried to align the cities' clocks by sending pulses of air in copper tubes beneath the streets. Many people resisted. New York City mayor Franklin Edson argued for the conventionality of time, pushing back against those who fulminated against the loss of Sun-borne time. In Paris, the poet Georges de Porto-Riche sued the city because he believed the mechanical pulses had destroyed his creativity. The opera includes a filmed melodrama about Martial Bourdin, a French anarchist who, in 1894, tried to blow up the Royal Observatory in Greenwich, London, where the global standard time was kept.

How does Einstein fit in?

The fury was even greater when the young Einstein proposed the radical idea that time is not absolute. His 1905 special theory of relativity implied that every person in motion carries a private time. This is encapsulated in the twin paradox that quickly emerged from Einstein's theory, in which a space traveller returns to Earth after a high-speed journey to find that his stay-at-home twin is older than he is. In the opera, we allude to these imagined twins, as well as time slowed and accelerated, reversed and put back forward.

How does your forthcoming film *Containment* take on the topic of time?

Made with film-maker Robb Moss, *Containment* is a feature documentary about an extraordinary assignment demanded by the US government. To safely store nuclear waste underground, the Department of Energy has to decide how to warn future populations that the waste is there for a period not less than 10,000 years. That is a long span to consider: 10,000 years ago, prehistoric monument Stonehenge was science-fiction far in the future.

What did you learn from working with Kentridge?

One of the great pleasures of the collaboration has been starting with the science but following chains of association far beyond what I do as a physicist, historian or even film-maker. At one point, we began with a question that divides scientists today: if you throw an encyclopaedia into a black hole, is the information gone forever or does it somehow survive? One day, inspired by a player-piano in William's studio, we realized that we could project light through rolls of perforated paper to invoke information falling into a black hole. We also made a parade of silhouettes marching into darkness, including two men struggling over a clock. Time is about physics, of course, but it is also, even for the most hardboiled scientist, about mortality, and the unshakable anxiety about disappearing forever. ■

INTERVIEW BY JASCHA HOFFMAN

This interview has been edited for length and clarity.

Correspondence

Climate justice more vital than democracy

Decision-making based on social-justice principles could be more effective than democratic efforts against climate change (see N. Stehr *Nature* **525**, 449–450; 2015).

Democratic decision-making involves multiple stakeholders, and democracy emphasizes the mutual roles of actors: all preferences are treated as equal. In many regions of the world, however, the results of democratic choices can be strongly influenced by power relations and inequitable social arrangements, owing to differences in economic development, access to technology and knowledge.

Elites may use democratic processes to entrench their status or encroach on other social goals (B. Sovacool *et al. Nature Clim. Change* **5**, 616–618; 2015). This can lead to incremental or undesirable results, which might explain why large democratic nations such as the United States continue to oppose progressive climate legislation.

In our view, sound climate and energy planning should not treat all stakeholders in the same way. Instead, preferences and roles should be weighted to consider criteria related to equity, due process, ethics and other justice principles. This would ensure that stakeholder discussions and resulting policies serve to eradicate, rather than exacerbate, socio-economic vulnerability to a changing climate.

Jingzheng Ren, Michael Evan Goodsite *University of Southern Denmark, Odense, Denmark.*
Benjamin K. Sovacool *Aarhus University, Herning, Denmark.*
benjaminso@auhe.au.dk

Interdisciplinarity: two more principles

As scientists in a sustainable-development research institute, our experience has demonstrated two further principles that are crucial for successful

interdisciplinary research (see R. R. Brown *et al. Nature* **525**, 315–317; 2015).

First, solving real-world problems requires integration of the social and biophysical realms. Key to this is the explicit consideration of differences in scale. For example, water governance has a typical time horizon of years, yet hydrological processes operate on timescales from seconds to millennia.

Second, a critical mass of interdisciplinary scholars is necessary to become a credible counterpart for disciplinary researchers and to catalyse progress. Our institute employs 150 interdisciplinary researchers, who engage much greater numbers of single-discipline researchers in interdisciplinary projects.

Martin J. Wassen, Marko P. Hekkert *Utrecht University, the Netherlands.*
m.j.wassen@uu.nl

Interdisciplinarity: topping the charts

We offer ideas on why India ranks highest and Brazil fifth — above the United States and Europe — in terms of the numbers of interdisciplinary research papers that they publish in Elsevier journals (see *Nature* **525**, 306–307; 2015).

India's researchers are surrounded by dynamic ecological, social and economic problems. They therefore naturally turn to solution-oriented research that is informed by the complexity of their environment. Although interdisciplinary researchers still struggle for acceptance in India's traditional, rigidly structured university departments, new universities with centralized programmes are offering fresh and tempting perspectives.

In Brazil, the number of interdisciplinary graduate programmes is growing twice as fast as the number of disciplinary courses. The country's national accreditation system changed

15 years ago, when many of its master's and PhD programmes failed to fit any of their set 46 disciplinary slots.

Marcel Burszty *University of Brasília, Brazil.*
Seema Purushothaman Azim *Premji University, Bangalore, India.*
marcel@unb.br

Restore vaccine trust in Japan

Vaccination coverage for human papilloma virus (HPV) in Japan has plummeted following unconfirmed media reports of severe adverse reactions such as chronic pain and cognitive decline (see S. J. B. Hanley *et al. Lancet* **385**, 2571; 2015). In response, and presumably to avoid lawsuits, the government has suspended its proactive vaccination recommendation.

Several scandals in the past few decades have fuelled the Japanese public's scepticism. These include health hazards associated with some drugs and vaccines, widely publicized cases of misconduct among researchers, and questionable connections between medical professionals and pharmaceutical companies (see T. Tanimoto *et al. Nature* **512**, 371; 2014).

At the heart of the debate is Japan's compensation scheme for vaccine injury, which depends on a causal relationship being confirmed by the government. Claimants may also seek further compensation through private litigation against the government. In most countries, claimants must pursue one scheme or the other, not both.

In our view, the government should introduce a no-fault compensation system, following the lead of other countries, including the United States and Scandinavian nations (see C. Looker and H. Kelly *Bull. World Health Organ.* **89**, 371–378; 2011). Exemption from liability would make for a fairer and more transparent judgement in policymaking and

pave the way for a new Japanese vaccination programme in which public trust is restored.

Tetsuya Tanimoto, Eiji Kusumi *Navitas Clinic, Tokyo, Japan.*
Claire Leppold *University of Edinburgh, UK.*
tetanimot@yahoo.co.jp

Some rules for behavioural science

US President Barack Obama's executive order of 15 September, 'Using behavioral science insights to better serve the American people', is intended to strengthen the policymaking process. To prevent the Obama administration from falling victim to agendas backed by pseudoscience, I suggest that it should consider only studies that abide by certain established rules.

The rules were instituted in behavioural and experimental economics research in universities, largely in response to scientific misconduct by governments and academics. They stipulate that an independent ethics committee must review research involving human subjects before the start of any experiment; that participants should first give their informed consent; and that they should receive appropriate remuneration to ensure that their decisions reflect genuine preferences. Moreover, deception should never be used in experimental design or in the process of scientific enquiry.

In the interests of scientific transparency, studies that lead to policy change should be peer reviewed and experimental design details and data made publicly available.

Such guidelines would help to ensure that studies are relevant, rigorous, ethical and legal. They would also provide a means to replicate and test studies for bias, efficacy and accuracy.

Jason A. Aimone *Baylor University, Waco, Texas, USA.*
jason_aimone@baylor.edu

William E. Paul

(1936–2015)

A leading force in immunology.

William Erwin Paul was a major contributor to the development of modern immunology. He helped to transform cytokine biology, the study of small proteins involved in cell signalling, from crude assessments of uncharacterized cellular 'factors' into a science involving precise quantitative molecular analyses. He also elucidated the mechanisms controlling the production of antibodies — proteins that recognize and bind to specific antigens such as bacteria and viruses — and provided insights into how antigens are recognized by T cells, a type of white blood cell.

Paul, who died on 18 September, was born in 1936 in Brooklyn, New York. He prided himself on the fact that his higher education was at public rather than private institutions. He obtained his undergraduate degree in 1956 from Brooklyn College and his medical degree in 1960 from the State University of New York's Downstate Medical Center, also in Brooklyn. While at medical school, he married Marilyn Heller.

During a two-year medical residency at Massachusetts Memorial Hospitals (now Boston Medical Center), Paul worked on amyloidosis, a rare disease that occurs when a protein called amyloid builds up in tissues. This led to his first paper, in *Nature*. More than 600 publications would follow.

Paul joined the US National Institutes of Health (NIH) in 1962. While working in the endocrinology branch of the National Cancer Institute in Bethesda, Maryland, he helped to establish that a chemotherapy drug called methotrexate was extremely effective in treating women with choriocarcinoma, a cancer that usually occurs in the placenta. He also helped to develop radioimmunoassays for hormones — *in vitro* techniques used to measure hormone levels in the blood. These studies were just a prelude to the real blossoming of Paul's research career.

Paul had developed a strong interest in immunology as a student. He recounted in a memoir that he was smitten by "a slender volume of essays by Michael Heidelberger, the father of quantitative immunochemistry", which he had read while riding a tram in Brooklyn. Paul sought advice from colleagues on the best places to train in immunology, and in 1964 joined Baruj Benacerraf's group at New York University in New York City.



(Benacerraf later won the 1980 Nobel Prize in Physiology or Medicine for his work on the genetics of the immune response.)

Paul moved with Benacerraf to the NIH in 1968. When Benacerraf left to become chair of pathology at Harvard Medical School in Boston in 1970, he encouraged the appointment of Paul as his successor. Paul remained chief of the Laboratory of Immunology at the National Institute of Allergy and Infectious Diseases until his death.

Under his leadership, the lab emerged as a premier centre of immunology research. This was a result of both Paul's remarkable scientific accomplishments and his astute guidance of the department. He gave both support and independence to those he recruited to faculty positions (myself included).

Paul is best known for his research on the cytokine interleukin-4 (IL-4), which he discovered with Maureen Howard, a post-doctoral fellow in his laboratory. He went on to show how the cytokine both helped to mobilize the body's defence mechanisms against parasitic worms and played a part in the production of allergic symptoms. He also mapped out the IL-4-dependent signalling pathway, the genetic control of the cytokine's expression and the factors that prompted certain types of T cell to produce

it. This all-encompassing and rigorous approach became the prototype for studies on immune-system cytokines.

Paul's impact on science and health went further. As the first director of the NIH Office of AIDS Research from 1994 to 1997, he set the agenda for research on HIV/AIDS in the United States. As Mark Harrington, executive director of an advocacy group called the Treatment Action Group, noted, Paul's work led to "the emergence of highly active antiretroviral therapy, which is now being taken by over 15 million people around the world and is responsible for saving millions of lives and preventing millions of infections".

Paul headed up several societies and served on innumerable advisory panels, prize committees and editorial boards. For 31 years, he served as editor of the *Annual Review of Immunology*; for much of that time, it was the most cited publication in biomedicine. He was also widely known for his textbook *Fundamental Immunology*, first published in 1984.

Bill was highly respected for being both a scholar and a gentleman. He could find flaws in any seminar he attended, but when asking questions, he made the speaker feel as if they were being praised, not skewered.

He remained his optimistic, ebullient, engaged self to the end. Only weeks before his death, he was calling colleagues to discuss ongoing research; attending to his roles as mentor, colleague and administrator; and completing his book *Immunity* (Johns Hopkins Univ. Press; 2015), in which he explains to a lay audience how the same system that defends the body can also cause autoimmune diseases.

His intellectual strength, encyclopaedic knowledge and laser-like focus on science — rather than on scientific politics — were unique. Friends, colleagues and immunologists around the globe mourn his passing. ■

Ronald N. Germain is chief of the Laboratory of Systems Biology at the US National Institute of Allergy and Infectious Diseases, National Institutes of Health, Bethesda, Maryland, USA. He started working with Bill Paul in the Laboratory of Immunology in 1982, and remained a close colleague and friend for more than 30 years. e-mail: rgermain@niaid.nih.gov

RONALD N. GERMAIN

NEUROSCIENCE

Decrypting a brain enigma

The combined neuronal activity of two seemingly opposite types of Purkinje cell in the brain's cerebellum has been found to be required to control the jerky eye movements known as saccades in monkeys. [SEE LETTER P.439](#)

KAMRAN KHODAKHAH

The brain has tremendous information-processing power and computational capacity. Neuroscientists have made great efforts to unravel how the brain encodes, decodes and processes information, but deconstructing these computations is a difficult task. On page 439 of this issue, Herzfeld *et al.*¹ use monkeys to masterfully unearth how one of the brain's most intriguing structures, the cerebellum, efficiently encodes — and perhaps subsequently decodes — the information needed to control the quick, jerky eye movements known as saccades, which occur as the eye explores a scene.

Deciphering how the brain computes is exceptionally difficult, because it requires a reverse-engineering approach. Without any prior knowledge of how the brain's circuits are designed, neuroscientists need to dissect these circuits to understand their function and computational principles. Consider the heroic efforts needed to reverse-engineer the Enigma machine, which was used during the Second World War to encrypt and decrypt military messages. By analogy, each of the brain's different computational units (neurons or neuronal circuits that process the information contained in their inputs to determine appropriate outputs) can be considered to be like an Enigma machine, with its own algorithm and code. Although Enigma's purpose was to encrypt and subsequently decrypt the same message, brain circuits process the information that they receive and often completely transform it to generate a new message. Furthermore, we do not fully understand the different types of information processing that occur in different brain regions.

However, this problem can be made more manageable by the thoughtful selection of specific brain regions to study. The more information there is about the inputs and outputs of a brain region, the details of its neuronal circuitry, the nature of the information it processes, and how and why it transforms the information, the better the chance of cracking its codes and understanding the purpose of its computations. Fortunately, because different brain regions often use the same or similar principles, understanding the code that

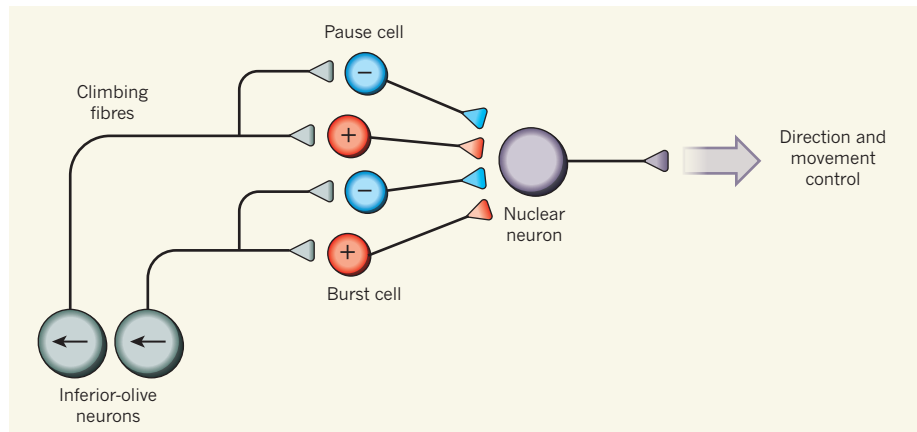


Figure 1 | Wiring up saccades. Neuronal projections called climbing fibres that originate in a brain region called the inferior olive encode the direction of saccades (quick, jerky eye movements that occur when exploring a scene; directional tuning indicated by arrows). Climbing fibres send input to Purkinje cells in the brain's cerebellum (although note that the activity of the Purkinje cells is primarily driven by projections called mossy fibres, which are not shown). The activity of one type of Purkinje cell, called burst cells, transiently increases with saccades (indicated by a + symbol), whereas the activity of another type, pause cells, ceases when saccades begin (indicated by a – symbol). Herzfeld *et al.*¹ suggest that, in monkeys, pause- and burst-cell signals are integrated by the nuclear neurons to which they signal to coordinate the movement of the eye during saccades. In addition, the authors used computer simulations to show that Purkinje cells that have climbing-fibre inputs with the same saccade directional tuning may target the same nuclear neuron. In this way, a simple neuronal circuit controls both the direction and movement of the eye during saccades. (Figure adapted from ref. 1.)

underlies one brain circuit can shed light on similar computations in other regions.

One region that has great potential for helping us to understand the brain is the cerebellum, which has a fairly simple anatomy. The cerebellum is mainly composed of repeats of the same computational circuit, at the core of which is a type of neuron called a Purkinje cell. Among other roles, the cerebellum controls and coordinates movements, and the anatomical connections involved in these functions are reasonably well delineated. These features, along with the ease with which the cerebellum's motor outputs can be quantitatively monitored, make it an ideal structure for deciphering the computational principles of some of the brain's Enigma machines (neuronal circuits).

The cerebellum is essential for the accurate control of saccades², the jerky eye movements that endow us with high-resolution vision by ensuring that the high-acuity region of the retina is exposed to the most important components of an image. Electrical recordings from Purkinje cells have long since revealed

that the activity of some of these neurons changes with saccades^{3–5}. But the individual activity of each saccade-related Purkinje cell is a poor predictor of saccade kinematics^{3–6} — the moment-to-moment speed of the eye during saccades. This is in stark contrast to the cerebellar control of smooth eye pursuit, in which the activity of individual Purkinje cells reliably predicts, in real time, the movement of the eye as it follows a moving object^{7,8}.

The failure of individual Purkinje cells to accurately predict saccade kinematics has led to speculation that the combined population-wide activity of Purkinje cells might more effectively encode the required information^{6,9}. Two types of saccade-related Purkinje cell have been identified: burst cells, whose activity transiently increases with saccades, and pause cells, whose spontaneous pre-saccade activity ceases when saccades begin. Although the activity of pause cells has been largely ignored, the population-wide activity of burst cells has been examined and found not to be predictive of saccade kinematics⁶. By analysing the

activity of Purkinje cells in monkeys as the animals made saccades, Herzfeld *et al.* made a breakthrough. They found that although neither the pause cells nor the burst cells predicted saccades as an individual population, their combined activity accurately predicted saccadic eye movements.

How might the brain decode the information encoded by the activity of the Purkinje-cell population? Herzfeld and colleagues suggest a plausible mechanism based on the anatomical organization of the cerebellum. Most of the cerebellar output is from clusters of neurons compacted into structures called cerebellar nuclei. Each nuclear neuron receives converging information from about 50 Purkinje cells. This organization⁹, combined with the speed at which Purkinje cells transfer information¹⁰, might enable nuclear neurons to faithfully integrate the population-wide activity of Purkinje cells in real time while preserving its temporal profile (Fig. 1). This possible decoding mechanism warrants careful experimental and theoretical scrutiny.

How the cerebellum encodes the direction of saccades is another long-standing puzzle, and Herzfeld and colleagues again offer a solution. Each Purkinje cell receives an input from

a climbing fibre — the output projection of neurons in the inferior olive, a brain region outside the cerebellum. The activity of climbing fibres is tuned to the direction of saccades¹¹. The authors used computer simulations to explore the possibility that climbing fibres dictate the functional organization of the cerebellar circuitry. They postulated that Purkinje cells that receive inputs from climbing fibres with the same directional tuning converge on the same nuclear neurons (Fig. 1). Remarkably, they found that when this organization is modelled, the activity of the Purkinje-cell population encodes both the real-time motion of the eye and the direction of the saccade, through a gain field — a multiplicative encoding mechanism found in the brain's cortex. Although purely speculative, this is a graceful solution to a difficult puzzle and is well worth experimental validation.

In addition to its role in motor control, the cerebellum has also been implicated in cognitive tasks. It is likely that many of the computational principles that the cerebellum uses for motor coordination are also implemented in its cognitive functions. It is equally likely that there are more, perhaps specialized, algorithms dedicated to its non-motor tasks.

Some of the brain's toughest and most elegant Enigma machines remain to be cracked, and the cerebellum might offer unique advantages for tackling them. ■

Kamran Khodakhah is at the Dominik P. Purpura Department of Neuroscience, Albert Einstein College of Medicine, New York, New York 10461, USA.
e-mail: k.khodakhah@einstein.yu.edu

1. Herzfeld, D. J., Kojima, Y., Soetedjo, R. & Shadmehr, R. *Nature* **526**, 439–442 (2015).
2. Barash, S. *et al. J. Neurosci.* **19**, 10931–10939 (1999).
3. Helmchen, C. & Büttner, U. *Exp. Brain Res.* **103**, 198–208 (1995).
4. Ohtsuka, K. & Noda, H. *J. Neurophysiol.* **74**, 1828–1840 (1995).
5. Kase, M., Miller, D. C. & Noda, H. *J. Physiol. (Lond.)* **300**, 539–555 (1980).
6. Thier, P., Dicke, P. W., Haas, R. & Barash, S. *Nature* **405**, 72–76 (2000).
7. Shidara, M., Kawano, K., Gomi, H. & Kawato, M. *Nature* **365**, 50–52 (1993).
8. Medina, J. F. & Lisberger, S. G. *J. Neurosci.* **27**, 6832–6842 (2007).
9. Gad, Y. P. & Anastasio, T. J. *Neural Networks* **23**, 789–804 (2010).
10. Person, A. L. & Raman, I. M. *Nature* **481**, 502–505 (2012).
11. Soetedjo, R., Kojima, Y. & Fuchs, A. F. *J. Neurophysiol.* **100**, 1949–1966 (2008).

CLIMATE SCIENCE

The long future of Antarctic melting

Simulations show that melting of the Antarctic ice sheet in response to climate change could raise the global sea level by up to 3 metres by the year 2300 and continue for thousands of years thereafter. [SEE LETTER P.421](#)

ALEXANDER ROBEL

Most projections of Antarctic ice melting in response to climate change extend a maximum of a few centuries into the future, a timescale that has clear relevance to immediate human affairs. But to capture the total Antarctic contribution to

sea-level rise caused by climate change, it is necessary to consider the possibility that ice-sheet mass loss will continue for thousands of years. On page 421 of this issue, Golledge and colleagues¹ present multi-millennial ice-sheet simulations in which Antarctica continues to contribute significantly to global mean sea-level rise for more than 1,000 years, long after ocean

and air temperatures have stopped increasing. The authors' simulations also show that ice-shelf melting driven by ocean warming over the next 100–300 years is a critical factor in determining the total future rise in global sea level.

The spectre of the imminent and rapid loss of ice from the West Antarctic Ice Sheet has been raised by dramatic events such as the collapse of the Larsen B ice shelf in 2002 (Fig. 1). Nevertheless, the contribution of the Antarctic ice sheet to global sea-level rise is currently small in comparison with that of other sources, although it is increasing at an accelerating rate².

The glaciologist John Mercer was the first to suggest that past recessions of the West Antarctic Ice Sheet were driven by fluctuations in climate that increased air temperatures over ice shelves to above freezing point³. He proposed that ice-shelf melting causes an increase in the flow of ice from land that occurs through other processes. In many coastal

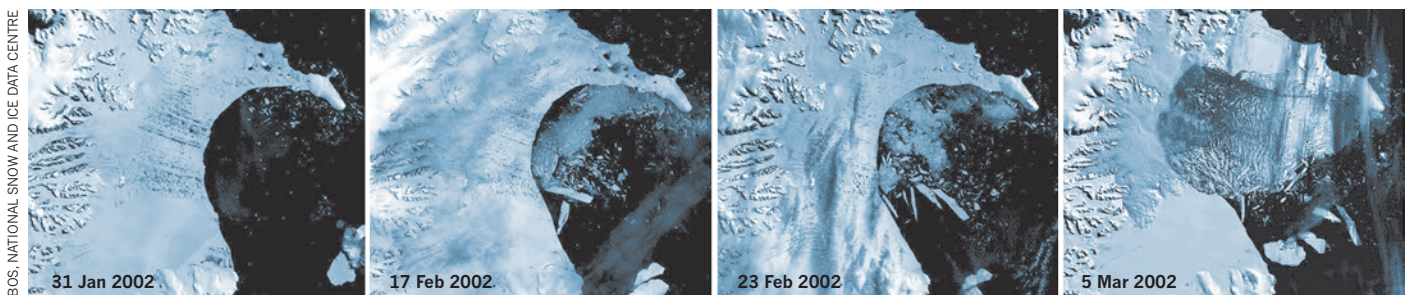


Figure 1 | Ice-shelf collapse. These satellite images show the rapid break-up of 3,250 square kilometres of the Larsen B ice shelf in Antarctica in 2002. Golledge *et al.*¹ have simulated the retreat of the Antarctic ice sheet in response to climate change over the next few millennia.

regions of Antarctica, ice shelves are laterally constrained by bays. Confirming Mercer's suspicions, numerical models^{4,5} have since shown that contact between the ice shelf and bay sidewalls exerts a restraining force on the flow of ice from land, leading to slower flow and reduced ice calving into the ocean. When the climate warms, ice shelves may melt, reducing the restraining contact with bay sidewalls, accelerating ice flow from land and enhancing mass loss through ice calving.

Estimating the likelihood that climate change will cause widespread Antarctic ice-shelf loss, and the duration of ice-flow acceleration in response to such a loss, is crucial for constraining future sea-level rise. However, accurately modelling ice sheets, ice shelves and their interactions with the ocean and atmosphere is computationally intensive. Consequently, most sea-level projections from models that explicitly include ice-sheet flow and ice-shelf melting extend only a few centuries into the future, or include only part of the Antarctic ice sheet. Existing multi-millennial projections of the ice sheet's contribution to sea-level rise are based on relatively simple statistical relationships between global temperature and sea level derived from palaeoclimatic and instrumental data^{6,7}.

In contrast to earlier studies, Golledge and colleagues use a comprehensive ice-sheet model, with forcing of precipitation, ocean and air temperature from global climate models, to simulate mass loss from the Antarctic ice sheet from the present to the year 5000 under various scenarios of climate change. They find that ocean warming, rather than atmospheric warming or changes in precipitation, is the dominant driver of mass loss. Ocean warming causes loss of large Antarctic ice shelves and a sustained acceleration in ice-sheet discharge into the ocean in all of the scenarios, except for one that includes the most extreme reduction in greenhouse-gas emissions compared with 1990 emission levels. In the scenarios in which ice shelves are lost, the long-term contribution of the Antarctic ice sheet to global sea-level rise ranges from roughly 3 to 9 metres. The majority of that contribution comes after 2300, with enhanced rates of sea-level rise lasting until at least 3000 — long after ocean temperatures have stabilized.

A previous theoretical analysis⁸ found that high model resolution is needed around the grounding line — which marks the transition from a grounded ice sheet to a floating ice shelf — to simulate rapid grounding-line migration accurately. Golledge and co-workers instead used a fairly coarse model resolution near the grounding line, because of the computational constraints of performing multi-millennial simulations of the entire Antarctic ice sheet.

For each climate-change scenario, they simulated an upper bound of sea-level rise using a scheme for correcting errors that arise

as a result of the coarse resolution near the grounding line, and a lower bound that was calculated by not applying such corrections. The sensitive dependence of projections of sea-level rise on such correction schemes is vividly demonstrated by the gap of several metres between the simulated lower and upper bounds. Efforts to develop more-accurate ways of representing the grounding line⁴ that can be incorporated into multi-millennial ice-sheet simulations should continue, and will lead to a considerable reduction in uncertainties in long-term sea-level projections.

Golledge *et al.* ultimately confirm the suspicions of earlier glaciologists that the fate of ice shelves largely determines whether Antarctica contributes less than 1 metre or up to 9 metres to long-term sea-level rise. Although ocean warming is responsible for most ice-shelf melting in these simulations, other studies^{9,10} have suggested that warm air temperatures lead to water ponding on the surface of ice shelves. This eventually causes hydrofracture (deepening of fractures by the drainage of meltwater ponds) and the rapid break-up of ice shelves. Capturing such complex processes in models is difficult, but, in one study that included these effects¹¹, Antarctic ice-shelf loss is even

more rapid and widespread than in the current simulations. If such a rapid ice-shelf break-up does occur, then Golledge and colleagues' simulations might represent a best-case scenario for future sea-level rise. ■

Alexander Robel is in the Division of Geological and Planetary Sciences, California Institute of Technology, Pasadena, California 91125, USA, and in the Department of Geophysical Sciences, University of Chicago. e-mail: robela@caltech.edu

1. Golledge, N. R. *et al.* *Nature* **526**, 421–425 (2015).
2. Harig, C. & Simons, F. J. *Earth Planet. Sci. Lett.* **415**, 134–141 (2015).
3. Mercer, J. H. *IAHS Publ.* **79**, 217–225 (1968).
4. Golberg, D., Holland, D. M. & Schoof, C. J. *Geophys. Res.* **114**, F04026 (2009).
5. Gudmundsson, G. H., Krug, J., Durand, G., Favier, L. & Gagliardini, O. *Cryosphere* **6**, 1497–1505 (2012).
6. Schaeffer, M., Hare, W., Rahmstorf, S. & Vermeer, M. *Nature Clim. Change* **2**, 867–870 (2012).
7. Levermann, A. *et al.* *Proc. Natl Acad. Sci. USA* **110**, 13745–13750 (2013).
8. Schoof, C. J. *Fluid Mech.* **573**, 27–55 (2007).
9. Scambos, T. A., Hulbe, C. & Fahnestock, M. A. *Antarctic Res. Ser.* **79**, 79–92 (2003).
10. Banwell, A. F., MacAyeal, D. R. & Sergienko, O. V. *Geophys. Res. Lett.* **40**, 5872–5876 (2013).
11. Pollard, D., DeConto, R. M. & Alley, R. B. *Earth Planet. Sci. Lett.* **412**, 112–121 (2015).

IMMUNOLOGY

A bacterial nudge to T-cell function

The epithelial cells that line the intestine have been found to sense tight attachment of bacteria, and to respond by producing proteins that shape the effector functions of the immune system's T_H17 cells.

SHAI BEL & LORA V. HOOPER

Our intestines contain trillions of bacteria that support health by promoting digestion, blocking invading microorganisms and synthesizing certain vitamins. Although the epithelial cells that line the gut separate these bacteria from deeper tissues, resident bacteria influence the development of immune cells that reside beneath the epithelium. This presents a puzzle: how do bacteria that are confined to the gut lumen communicate their presence to these immune cells? Writing in *Cell*, Atarashi *et al.*¹ and Sano *et al.*² show that the gut epithelial lining provides a conduit for this communication by sensing bacterial attachment and producing proteins that guide immune-cell development.

Among the hundreds of bacterial species that live in the intestine, a select few are especially adept at stimulating the immune system. One such bacterial group is the segmented

filamentous bacteria (SFB). SFB have a distinctive appearance, being composed of long, segmented chains, or filaments. The filament ends are buried in the outward-facing membranes of epithelial cells and thus the bacteria are attached tightly to the intestinal surface³. Despite their menacing appearance under the microscope (Fig. 1a), the bacteria do not invade the epithelial barrier. But the tight attachment of SFB to the gut surface coincides with a remarkable ability to shape the development and function of subepithelial immune cells. These include T_H17 cells, a key immune-cell group that is the focus of the present studies.

Numerous cell types promote immunity to infection. 'Helper' T cells (T_H cells), characterized by the presence of the CD4 protein on their surface, regulate immunity by producing cellular signalling molecules called cytokines. T_H cells can differentiate into several lineages that differ in the cytokines they make and the

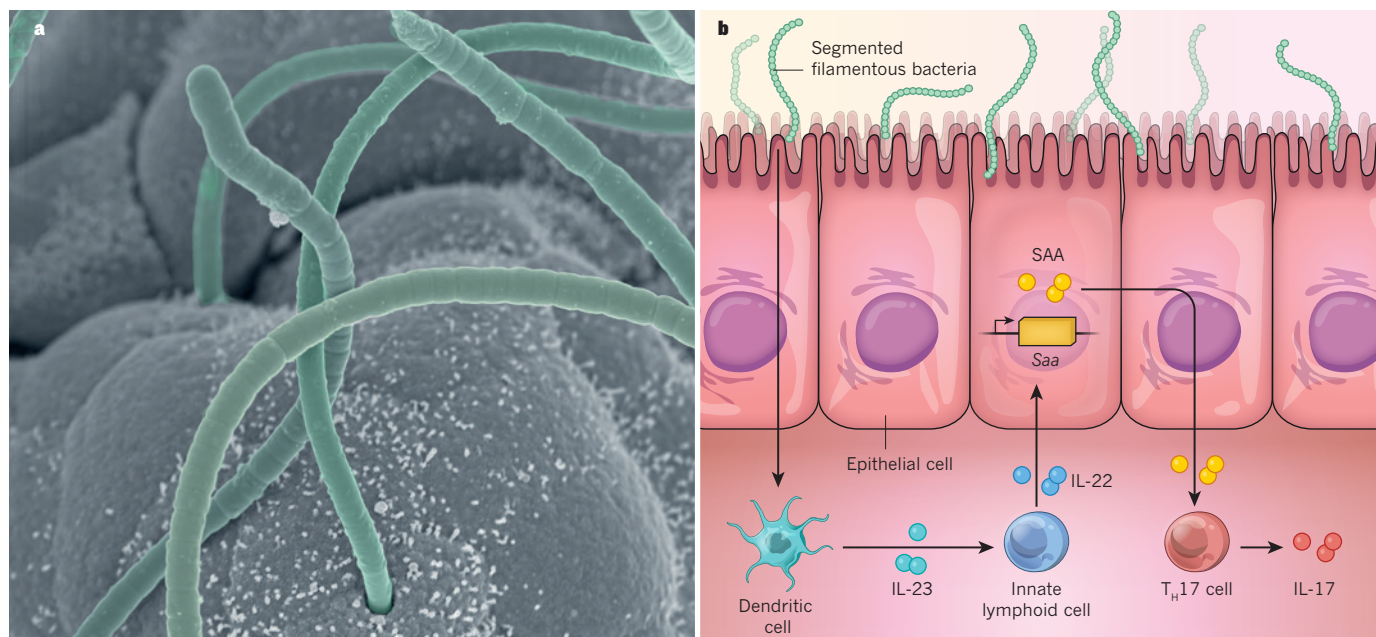


Figure 1 | Bacterial attachment to the intestinal surface shapes T_H17 -cell function. **a**, Segmented filamentous bacteria (SFB; coloured green) attach tightly to the epithelial cells that line the outer surface of the intestine in mice. **b**, Atarashi *et al.*¹ provide evidence that this tight attachment of SFB to the intestinal epithelium shapes the function of T_H17 cells of the immune system by inducing epithelial-cell expression of the protein serum amyloid A (SAA).

Sano *et al.*² show that SAA spurs production of the cell-signalling molecule IL-17 by differentiated T_H17 cells. The authors also unravel a regulatory circuit, involving other subepithelial immune cells, that governs epithelial *Saa* expression in this context. The circuit includes IL-23 (probably produced by dendritic cells in response to bacterial signals), which stimulates innate lymphoid cells to produce IL-12. IL-22 then induces *Saa* expression in epithelial cells.

infections that they target. T_H17 cells produce the cytokines IL-17 and IL-22, which protect against extracellular pathogens and fungi. But T_H17 cells also have a dark side — too many of them in the wrong place can provoke autoimmune and inflammatory diseases such as rheumatoid arthritis, multiple sclerosis and Crohn's disease⁴.

Working together, the two research groups behind the new papers previously discovered that SFB promote a marked accumulation of IL-17-secreting T_H17 cells in the intestinal tissues⁵. This accumulation correlated with increased resistance to bacterial infection as well as a heightened susceptibility to autoimmune disease^{5,6}. But the mystery remained of how a non-invasive bacterial species alters T-cell development on the other side of the epithelial barrier.

Atarashi *et al.* provide compelling evidence that the physical interaction of SFB with the gut epithelium is essential to spur T_H17 -cell development. The authors first noticed that SFB have a remarkable fidelity towards their animal host. When SFB from either mice or rats were introduced into microbiologically sterile (germ-free) mice, only the mouse SFB attached to the gut surface and induced T_H17 -cell accumulation. But when the same bacteria were introduced into germ-free rats, it was the rat SFB that adhered tightly to the host and induced a T_H17 response. This suggested that adhesion of SFB to intestinal epithelial cells is host-specific, and that it correlates with T_H17 -cell accumulation.

To nail down whether this was a causal

relationship, the authors studied two other organisms that attach to the epithelium and induce T_H17 -cell accumulation: *Citrobacter rodentium* and *Escherichia coli*. Unlike SFB, the genomes of these bacteria can be altered in the laboratory, and so Atarashi and colleagues deleted the bacterial genes that are required for epithelial attachment. They found that there were fewer T_H17 cells in the guts of mice colonized with the defective bacteria, suggesting that bacterial attachment causes T_H17 -cell accumulation.

So how does bacterial attachment trigger the development of cells that lurk in subepithelial tissues? Both Atarashi *et al.* and Sano *et al.* analysed epithelial cells with and without attached SFB to see how they differed. Among other things, they found that epithelial cells with attached SFB made more of a protein called serum amyloid A (SAA), expression of which is induced during infection or injury.

The researchers had previously noticed that SFB induces SAA expression, and experiments with cultured T cells had suggested that SAA might coax T cells to differentiate into T_H17 cells⁵. In their present work, Sano *et al.* extend this observation to show that SAA actually induces committed T_H17 cells to produce more IL-17, indicating that SAA modifies the effector functions of these cells rather than their commitment to a particular lineage. The authors support this idea by studying mice engineered to lack SAA — although these mice have normal numbers of committed T_H17 cells, the cells make less IL-17 than cells from wild-type mice.

Sano *et al.* also unravel details about the molecular circuitry that regulates SAA production. They show that SFB bolster SAA expression through a roundabout route involving innate lymphoid cells (ILCs), a group of immune cells that are related to T cells but distinct in that they lack receptors for specific antigen molecules. The authors found that SFB prompt ILCs to produce IL-22, which in turn elicits SAA production by gut epithelial cells (Fig. 1b). Interestingly, a similar regulatory circuit governs the expression of key antimicrobial proteins that are lobbed like hand grenades from the outward-facing epithelial surface to minimize the attachment of both foreign and resident bacteria⁷. Sano and colleagues' findings suggest that this circuit not only regulates the response to bacterial attachment, but also readies subepithelial immune cells to deal with potential invaders.

The findings raise several fascinating questions. For example, how do epithelial cells sense a physical interaction with bacteria? Atarashi *et al.* provide tantalizing preliminary evidence that this might involve SFB triggering rearrangements of actin (a structural protein) that alter transcription of the *Saa* gene. And how does SAA activate IL-17 expression? SAA may act on T_H17 cells through a cytokine-like mechanism. An alternative possibility stems from the finding that SAA binds to retinol (vitamin A)⁸, which suggests that SAA could influence immune-cell development by delivering retinol to the cells.

A key motivation behind research on gut bacteria is to harness their functional properties

to improve human health. The present studies suggest the possibility of using tightly adherent but non-invasive bacteria to deliberately nudge the human intestinal immune system towards a more activated state. Such an approach could be particularly useful in preventing or treating infectious disease. Although SFB do not colonize humans, Atarashi *et al.* have taken a step in this therapeutic direction by identifying a group of 20 human-gut bacteria that, like

SFB, adhere to the intestinal surface and induce T_H17-cell accumulation in mice. These findings should inspire further investigation into how adherent non-invasive bacteria influence the human immune system. ■

Shai Bel and Lora V. Hooper are in the Department of Immunology, University of Texas Southwestern Medical Center, Dallas, Texas 75390, USA. L.V.H. is also at the

Howard Hughes Medical Institute.

e-mail: lora.hooper@utsouthwestern.edu

1. Atarashi, K. *et al.* *Cell* **163**, 367–380 (2015).
2. Sano, T. *et al.* *Cell* **163**, 381–393 (2015).
3. Davis, D. P. & Savage, D. C. *Infect. Immun.* **10**, 948–956 (1974).
4. Burkett, P. R., Meyer zu Horste, G. & Kuchroo, V. K. *J. Clin. Invest.* **125**, 2211–2219 (2015).
5. Ivanov, I. I. *et al.* *Cell* **139**, 485–498 (2009).
6. Wu, H. J. *et al.* *Immunity* **32**, 815–827 (2010).
7. Sanos, S. L. *et al.* *Nature Immunol.* **10**, 83–91 (2009).
8. Derebe, M. G. *et al.* *eLife* **3**, e03206 (2014).

emission is forbidden by angular-momentum-conservation arguments. This turned out to have astrophysical consequences: the two-photon decay of hydrogen and hydrogen-like atomic systems contributes to the strength of the continuum spectra of planetary nebulae at ultraviolet wavelengths⁷. The reverse process — the simultaneous absorption of two photons during the excitation of an atom or a molecule — has even broader impact: it forms the basis of two-photon laser-scanning fluorescence microscopy⁸, which is used in the biomedical fields and in tissue engineering⁹.

Now, some 30 years after the pioneering measurements of $\gamma\gamma$ -decay in transitions of ¹⁶O, ⁴⁰Ca and ⁹⁰Zr nuclei, Walz and co-workers have detected, for the first time, the $\gamma\gamma$ -decay of a nuclear quantum transition that competes with the allowed single γ -decay. The authors faced two experimental challenges: first, the former process is at least 100,000 times less probable than the latter one. Second, the γ -ray emitted in single γ -decay easily scatters off electrons in the detector material (a process known as Compton scattering), so that some of its energy is deposited in one detector and the rest in another, in a way that mimics the signature of the elusive $\gamma\gamma$ emission.

In an experimental tour de force, Walz *et al.* used an established γ -ray calibration standard, the ¹³⁷Ba nucleus in an excited state that is created when a caesium nucleus (¹³⁷Cs) undergoes a process known as β -decay. When the excited state of ¹³⁷Ba decays, it predominantly emits a 662-kiloelectronvolt γ -ray. To catch the rare, simultaneous $\gamma\gamma$ events, the authors surrounded the γ -ray source with five scintillation detectors (see Fig. 1a of the paper¹) that could register γ -rays at a time resolution better than 1 nanosecond. This fine resolution crucially helped the authors to reject the impostor, scattered single- γ emissions, which have different flight times from the true simultaneous $\gamma\gamma$ events. Furthermore, the authors used massive lead shielding to minimize the scattering of γ -rays from one detector into another.

After more than 50 days of measurements, Walz *et al.* reported a clear signal at the expected total energy (662 keV). It was registered simultaneously by detector pairs from various angles, thus yielding the first positive signal of $\gamma\gamma$ -decay in competition with single- γ

NUCLEAR PHYSICS

Sometimes γ -rays come in twos

Breakthrough measurements of a rare decay process in an excited barium nucleus pave the way for the development of techniques that probe the structure and decay modes of atomic nuclei. SEE LETTER P.406

ALEXANDRA GADE

The excited quantum states of atomic nuclei predominantly decay through the emission of short-wavelength, energetic electromagnetic radiation known as γ -rays — an excited nucleus emits a single γ -ray photon to reach a lower energy state. The energy radiated is equal to the difference in the energies of the initial and final quantum states. But on page 406 of this issue, Walz *et al.*¹ report the first observation of a rare decay mode, for an excited quantum state of a barium nucleus (¹³⁷Ba), that proceeds through the simultaneous emission of two γ -ray photons. Because this quantum state of ¹³⁷Ba usually decays through the much more common emission of a single γ -ray photon, Walz and colleagues' laboratory measurements represent a formidable exercise in finding a needle in a haystack, and have relevance to the theories of quantum electrodynamics and nuclear structure.

The $\gamma\gamma$ -decay was first described by the physicist Maria Göppert-Mayer on the cusp of the 1930s. In her doctoral thesis, she predicted^{2,3} that this rare electromagnetic decay mode is allowed in the framework of quantum electrodynamics: the energy of each of the two emitted γ -ray photons can take a value from a continuous energy distribution, so that the total energy equals the energy liberated as the nucleus undergoes the quantum transition from the initial to a final state.

It was not until the mid-1980s that this rare process was verified experimentally for nuclei. These pioneering measurements involved observations of the $\gamma\gamma$ -decay modes of excited states of calcium (⁴⁰Ca; ref. 4), zirconium⁴ (⁹⁰Zr) and oxygen⁵ (¹⁶O) nuclei,

and required a clever trick. Both the ground and the first excited states of these nuclei have zero total angular momentum. The fact that a γ -ray photon (or any photon) cannot have zero angular momentum means that the decay of one of these excited states through single γ -ray emission is strictly forbidden because of the requirements of the conservation of angular momentum. This realization allowed researchers to eliminate the usually dominant single γ -decay mode from consideration, and left processes that are experimentally easy to identify (such as conversion electron emission

The authors determined that the odds of $\gamma\gamma$ -decay occurring were only 1 in 487,805.

and internal pair formation) as the only other competing decay modes^{4,5}.

Although $\gamma\gamma$ emission from nuclei has so far been a rare beast, the corresponding two-photon emission following electronic transitions in atoms was measured and studied several decades ago. This is because the atomic process has advantageous energetics: the lower-energy photons (such as X-rays) produced are more easily absorbed by one detector than are nuclear γ -rays, which more frequently scatter from one detector into another. Atomic physicists were also aided by the experimental tools to hand. For example, lasers can be used to 'pump' an atom into an excited atomic state that is a candidate for two-photon decay or that decays to such a state.

As early as the 1940s, the physicists Gregory Breit and Edward Teller conjectured that two-photon emission is the most probable radiative-decay mode of the metastable 2s_{1/2} state of atomic hydrogen⁶, for which single-photon

emission. The authors determined that the odds of the former decay occurring were only 1 in 487,805, underlining the accomplishment in isolating such a feeble signal from a multitude of much more probable processes.

Other groups have been trying for several years to detect and characterize the $\gamma\gamma$ -decay of ^{137}Ba through experimental^{10–12} and theoretical¹⁰ work, but Walz and co-workers have made the first breakthrough. One of the remaining challenges is to measure precisely the energy distributions of the two emitted γ -rays at various emission angles.

Although Walz and colleagues' measurement provides a rare confirmation of an elusive nuclear process, $\gamma\gamma$ emission itself could be used to probe the structure of atomic nuclei. The $\gamma\gamma$ -decay of a nucleus occurs through a multitude of virtual intermediate nuclear states. The corresponding transitions that occur as a nucleus passes from its initial state

to a virtual state, and from there to the final ground state, can be calculated by models that aim to describe the structure of nuclei. The characteristics of $\gamma\gamma$ emission can therefore provide a tracer of the combined probability that a nucleus can transition through the multitude of possible pathways — these can never be fully measured in experiments that attempt to excite all individual intermediate states.

The previously measured probabilities of $\gamma\gamma$ transition between two states that each have zero total angular momentum are related to gross properties of the nucleus, such as electric polarizability and magnetic susceptibility⁵. In the case of the more-complex decay process studied by Walz *et al.*, developing a comprehensive theoretical model is a challenge. Walz and colleagues' data may help by acting as a benchmark against which theoretical nuclear wavefunctions of ^{137}Ba can be compared. ■

Alexandra Gade is at the National Superconducting Cyclotron Laboratory, Michigan State University, East Lansing, Michigan 48824-1321, USA.
e-mail: gade@nsl.msu.edu

1. Walz, C. *et al.* *Nature* **526**, 406–409 (2015).
2. Göppert, M. *Naturwissenschaften* **17**, 932 (1929).
3. Göppert-Mayer, M. *Ann. Phys.* **9**, 273–294 (1931).
4. Schirmer, J. *et al.* *Phys. Rev. Lett.* **53**, 1897–1900 (1984).
5. Kramp, J. *et al.* *Nucl. Phys. A* **474**, 412–450 (1987).
6. Breit, G. & Teller, E. *Astrophys. J.* **91**, 215–238 (1940).
7. Spitzer, L. & Greenstein, J. L. *Astrophys. J.* **114**, 407–420 (1951).
8. Denk, W., Strickler, J. H. & Watts, W. W. *Science* **248**, 73–76 (1990).
9. Helmchen, F. & Denk, W. *Nature Methods* **2**, 932–940 (2005).
10. Millener, D. J., Sutter, R. J. & Alburger, D. E. *Bull. Am. Phys. Soc.* **56**, Abstr. BAPS.2011.DNP.CF.8 (2011).
11. Moran, K. *et al.* *Phys. Rev. C* **90**, 041303 (2014).
12. Lister, C. J. *et al.* *Bull. Am. Phys. Soc.* **58**, Abstr. BAPS.2013.DNP.CE.3 (2013).

COGNITIVE DISORDERS

Deep brain stimulation for Rett syndrome

Mutations in the gene *MECP2* cause an intellectual-disability disorder called Rett syndrome. In a mouse model, electrical stimulation of deep brain regions is found to ameliorate some of the features of the syndrome. SEE LETTER P.430

STUART R. COBB

Childhood intellectual-disability disorders typically have a genetic cause and are notoriously resistant to treatment. In particular, there are almost no therapies that can alleviate deficits in core cognitive features such as learning, memory, concentration and communication. So far, the search for treatments has been dominated by

drug-based¹ and, to a lesser extent, gene-based therapies (see ref. 2, for example). However, on page 430 of this issue, Hao *et al.*³ provide evidence that electrical stimulation of a deep brain region can effectively reverse impaired learning in a mouse model of Rett syndrome — one of the leading causes of intellectual disability in girls.

Deep brain stimulation (DBS) has been used to treat various brain disorders, including

Parkinson's disease, chronic pain, obsessive-compulsive disorder and clinical depression⁴. Its use is well established in movement disorders that are otherwise resistant to treatment, but it remains highly experimental for most other neurological and psychiatric conditions, and has rarely been used in children. Using DBS entails implanting thin wire electrodes in the brain, to a depth and location determined by the symptoms being treated. An electrical device then delivers controlled impulses to stimulate local brain activity in a regulated fashion. However, the precise mechanism by which DBS works remains largely unknown.

Experimental evidence suggests that DBS may improve cognition in adult rats with memory impairment⁵, and can slow the rate of cognitive decline in people with Alzheimer's disease⁶. But Hao and colleagues' report marks the first demonstration that DBS might have the potential to tackle a childhood intellectual-disability disorder. The authors focused on Rett syndrome, which lies towards the more severe end of the intellectual-disability spectrum and

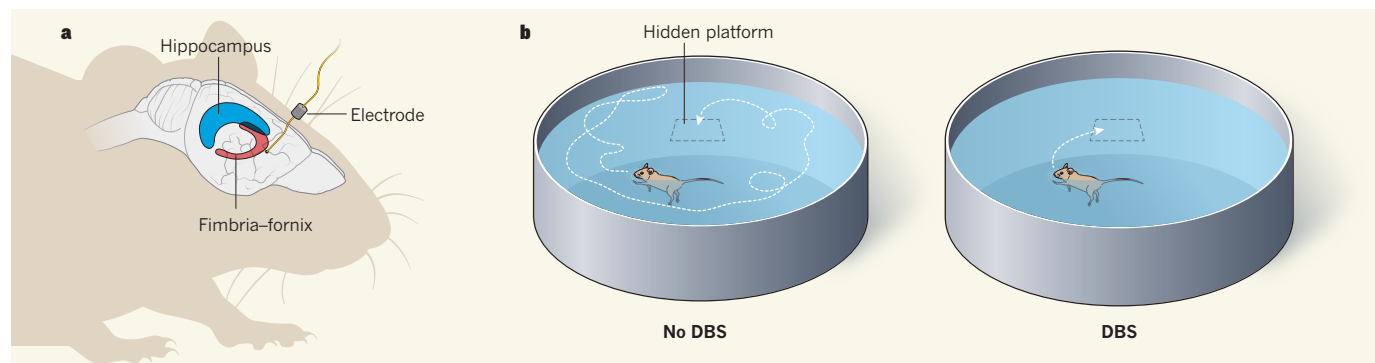


Figure 1 | Stimulating learning. **a**, Hao *et al.*³ treated mice that model Rett syndrome — an intellectual-disability disorder — with a procedure known as deep brain stimulation (DBS). The authors implanted electrodes into a nerve-fibre bundle in the brain called the fimbria-fornix, which connects the hippocampal regions responsible for some forms of

memory, both with each other and with deeper brain structures. They then stimulated this region electrically. **b**, Using DBS enhanced hippocampus-dependent learning and memory in Rett-syndrome mice, enabling treated mice to find a hidden platform more efficiently than those that did not receive DBS.

is almost always caused by mutations in the *MECP2* gene⁷. Mice that have mutations in this gene show intellectual disability and mimic other key features of Rett syndrome, including impairment of motor skills and breathing control, and other neurological effects.

Hao *et al.* targeted DBS to the fimbria-fornix (Fig. 1), an intersection of nerve fibres that connects the hippocampus regions in each brain hemisphere (which are involved in certain forms of memory), both with each other and with other brain structures. The authors stimulated the fimbria-fornix of mice for two weeks and performed sensitive behavioural tests three weeks later. In normal mice that received DBS, there was a modest improvement compared with unstimulated mice in forms of learning and memory that are dependent on hippocampal function. Crucially, however, Rett-syndrome mice showed a dramatic improvement in hippocampus-dependent spatial learning and contextual fear memory. For both of these cognitive functions, behavioural performance in Rett-syndrome mice was restored to levels indistinguishable from those seen in normal, unstimulated animals.

These exciting findings could have major implications, but caution is required when extrapolating from cognitive tests in mice to the treatment of humans with intellectual disability. The tests used in animal models tend to be very specific, and it is unclear to what extent cognitive functions measured in rodents, such as ability to perform in a maze test, correspond to the cognitive domains that need treating in people. Moreover, behavioural tests in rodents can be poor predictors of human responses. This is particularly true for tests of drugs that enhance cognition — there is a high attrition rate for these drugs, because many that have shown promise in cognitive tests in animals lacked efficacy in human clinical trials.

The authors' results indicate that the effects of fimbria-fornix DBS are restricted to specific cognitive domains, with no improvement in other Rett-syndrome-like symptoms, such as altered anxiety, pain sensitivity or motor control. This is perhaps not surprising, given the brain region targeted by Hao and colleagues. However, DBS has been effective in lessening these other defects, especially in motor disorders. So it is possible that altering the position of electrodes may ameliorate other symptoms of Rett syndrome.

The next steps are to analyse the mechanism by which DBS improves learning in Rett-syndrome mice, and to determine whether similar effects can be expected in people with Rett syndrome or other forms of intellectual disability. The main way in which memory traces are thought to be encoded and stored is through changes in the strength of the synaptic connections between neurons, and disruption of this synaptic plasticity is a hallmark of many neurodevelopmental disorders. Hao *et al.* show that long-term potentiation — a particular

type of synaptic plasticity that is required for certain forms of learning — is impaired in the hippocampus of Rett-syndrome mice but is boosted by DBS. Furthermore, DBS also enhanced the generation of neurons from stem cells that are resident in the hippocampus. Whether these effects are simply biomarkers of stimulation in hippocampal neural circuits, or whether they truly contribute to the procognitive action of DBS, awaits further experimentation.

Hao and colleagues' study is meaningful because finding effective treatments for the core symptoms of childhood intellectual-disability disorders is one of the great unmet medical challenges in contemporary neuroscience. Genetic studies indicate that a multitude of molecular abnormalities can cause intellectual disability^{1,8}, which probably explains the prevalence of these disorders and makes finding treatments so difficult — it is hard to envisage a drug therapy that would be effective for many forms of intellectual disability.

Although DBS is invasive, it is considered to

be safe and controllable⁴. Whether or not DBS can be widely applied to intellectual-disability disorders, and whether it may one day constitute a common treatment option, remains to be seen. Nonetheless, the current study will doubtless provide impetus for future studies in this direction. ■

Stuart R. Cobb is at the Institute of Neuroscience and Psychology, College of Medical, Veterinary and Life Sciences, University of Glasgow, Glasgow G12 8QQ, UK. e-mail: stuart.cobb@glasgow.ac.uk

1. Ghosh, A., Michalon, A., Lindemann, L., Fontoura, P. & Santarelli, L. *Nature Rev. Drug Discov.* **12**, 777–790 (2013)
2. Garg, S. K. *et al. J. Neurosci.* **33**, 13612–13620 (2013).
3. Hao, S. *et al. Nature* **526**, 430–434 (2015).
4. Lyons, M. K. *Mayo Clinic Proc.* **86**, 662–672 (2011).
5. Shirvaskar, P. R., Rapp, P. R. & Shapiro, M. L. *Proc. Natl Acad. Sci. USA* **107**, 7054–7059 (2010).
6. Laxton, A. W. *et al. Ann. Neurol.* **68**, 521–534 (2010).
7. Lyst, M. J. & Bird, A. *Nature Rev. Genet.* **16**, 261–275 (2015).
8. Srivastava, A. K. & Schwartz, C. E. *Neurosci. Biobehav. Rev.* **46**, 161–174 (2014).

NANOPHYSICS

Microscopic friction emulators

Cold ions sliding across periodic energy-potential patterns formed by lasers have been used to elucidate the physics of dry friction between crystals. Experiments with no more than six ions suffice to explore a vast domain of frictional forces.

DAVIDE MANDELLI & ERIO TOSATTI

The study of friction is a centuries-old, but still vibrant, subject. The archetypal friction problem of dry sliding between solids concerns the force that resist the relative lateral motion of crystal surfaces in contact. In this case, the periodic spacing of the atoms in the crystals eliminates the complexity introduced by the ill-defined nature of ordinary, non-crystalline surfaces. Studies of nanometre-scale systems through tip-based instruments, non-equilibrium theory, computer simulations and, most recently, the use of artificial friction emulators have revived this topic¹. Writing in *Science* and in *Nature Physics*, respectively, Bylinskii *et al.*² and Gangloff *et al.*³ follow up earlier theoretical suggestions^{4–6}, and present friction emulators formed from trapped ions that slide over optical lattices — periodic energy potentials created by the interference of laser beams.

Although dry sliding is a well-defined problem, sufficient inroads have not been made into it. This is because suitable experimental systems are scarce, but also because the theory

of frictional dynamics is mostly restricted to computer simulations, which are vivid, but incomplete, representations of the complex phenomenon of friction between real solid interfaces. Much of scientists' understanding of friction relies on apparently trivial, one-dimensional periodic potential models, such as the Prandtl–Tomlinson (PT) single-slider⁷ and the Frenkel–Kontorova (FK) sliding-chain descriptions⁸.

The PT model illustrates the switch from stick-slip friction — for example, that arising from the intermittent motion of chalk on a blackboard — to smooth sliding as surface corrugation decreases. The same model also describes the passage from large friction at high sliding speeds to vanishingly small friction (thermolubricity) at very low speeds^{7,9}. In the special case of mismatched lattices, the FK model describes the transition between frictionless (superlubric) motion and pinned frictional sliding (a regime in which it takes a large force to dislodge the lattice and nudge it forward).

What these models have achieved is a description of how properties such as corrugation, temperature, velocity and lattice-matching

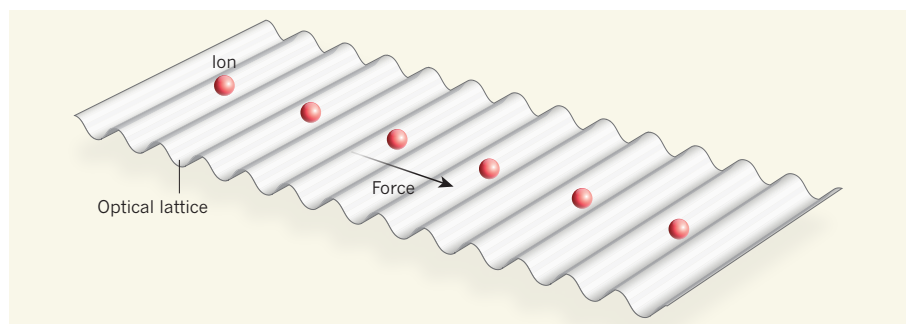


Figure 1 | Ion chain trapped in a lattice. Bylinskii *et al.*² and Gangloff *et al.*³ investigate the friction that arises when one or a few ions are forced to slide across optical lattices in which a corrugated energy potential is created by counter-propagating lasers. These simple experiments have provided benchmark tests of long-standing frictional models for sliding crystal surfaces. (Adapted from ref. 2.)

of materials in contact might influence and determine friction. In principle, the understanding gained from the models could allow researchers to control frictional forces, something desirable in many practical settings. However, despite the firm theoretical background^{7–9}, neither of these models has been tested experimentally. That is what emulators are meant to do.

The emulators of Bylinskii *et al.* and Gangloff *et al.* are based on short chains of trapped ionized atoms that, under the influence of an electric field, are forced to slide across a laser-generated optical lattice (Fig. 1). These techniques may seem arcane, yet they are accurate and powerful, because parameters such as the temperature, atom velocity and spacing, chain length and the amplitude of the lattice's potential can be flexibly adjusted across a vast range of values.

In this vein, Gangloff and colleagues present experiments involving one or two ions that slide across an optical lattice, a set-up that emulates the PT model to near perfection. In their single-ion set-up, the authors demonstrate that, even at microkelvin temperatures, there is a speed below which the friction between the sliding ion and the lattice vanishes — in agreement with thermodynamics. For higher speeds, stick-slip friction ensues and rises by more than a factor of 100 with increasing speed. No experiment involving real crystals can span this range of friction and speed. Although many of the authors' findings had been known from numerical simulations⁷, their emulator is superior to those — for example, it reproduces the theoretically expected dependence of friction on velocity¹⁰ much more accurately.

To adequately emulate the FK model, one would need an infinitely long ion chain instead of a single ion. All that Bylinskii *et al.* use are short chains of two to six ions sliding across an optical lattice — so is this just a baby step towards that goal? Not quite. By adjusting the distances between the ions, the authors tune the amount of mismatch between the chain and the corrugation of the lattice's periodic potential, and produce a dramatic

effect on the friction that they measure.

Starting from large friction for perfect chain–lattice matching, Bylinskii and colleagues observe a rapid decrease in friction as they increase a parameter that controls the amount of chain–lattice mismatch. This trend reflects the evolution from strong pinning friction towards lubricity or even superlubricity. Superlubricity, however, is reached only when the intensity of the lattice's potential falls below the value that demarcates the transition (known as the Aubry transition) between frictionless and pinned frictional sliding¹¹.

Studies of the Aubry transition are of interest, as has been suggested by theoretical studies of systems with long ion chains^{4,12}. It will be even more interesting to study this transition in short-chain systems such as those that Bylinskii *et al.* and Gangloff *et al.* use. Overall, it is surprising how much can be learnt about the physics of infinitely long systems from studies that involve just a few ions.

As always, experiments teach us more than

we anticipate. Clear-cut techniques, such as the cold-ion emulators reported in these two papers, provide insights into the complexity that underlies even the simplest act of friction involving a handful of ions. The physicist Philip Warren Anderson once said¹³, “more is different”. But in the case of the current papers, one could counter that dictum by saying that, sometimes, less can be different after all. ■

Davide Mandelli and Erio Tosatti are at the *International School for Advanced Studies (SISSA)*, 34136 Trieste, Italy. E.T. is also at the *Abdus Salam International Centre for Theoretical Physics, Trieste*, and the *Democritos National Laboratory, Istituto Officina dei Materiali, Consiglio Nazionale delle Ricerche, Trieste*.
e-mails: davide.mandelli@sisssa.it; tosatti@sisssa.it

1. Vanossi, A., Manini, N., Urbakh, M., Zapperi, S. & Tosatti, E. *Rev. Mod. Phys.* **85**, 529 (2013).
2. Bylinskii, A., Gangloff, D. & Vuletić, V. *Science* **348**, 1115–1118 (2015).
3. Gangloff, D., Bylinskii, A., Counts, I., Jhe, W. & Vuletić, V. *Nature Phys.* <http://dx.doi.org/10.1038/nphys3459> (2015).
4. Benassi, A., Vanossi, A. & Tosatti, E. *Nature Commun.* **2**, 236 (2011).
5. García-Mata, I., Zhurov, O. V. & Shepelyansky, D. L. *Eur. Phys. J. D* **41**, 325–330 (2007).
6. Pruttivarasin, T., Ramm, M., Talukdar, I., Kreuter, A. & Häfner, H. *N. J. Phys.* **13**, 075012 (2011).
7. Müser, M. H. *Phys. Rev. B* **84**, 125419 (2011).
8. Braun, O. M. & Kivshar, Y. *The Frenkel–Kontorova Model: Concepts, Methods and Applications* (Springer, 1998).
9. Krylov, S. Y. & Frenken, J. W. M. *Phys. Status Solidi B* **251**, 711–736 (2014).
10. Sang, Y., Dubé, M. & Grant, M. *Phys. Rev. Lett.* **87**, 174301 (2001).
11. Aubry, S. & Le Daeron, P. Y. *Physica D* **8**, 381–422 (1983).
12. Sharma, S. R., Bergersen, B. & Joos, B. *Phys. Rev. B* **29**, 6335 (1984).
13. Anderson, P. W. *Science* **177**, 393–396 (1972).

BEHAVIOURAL ECONOMICS

Visible inequality breeds more inequality

Experiments suggest that when people can see wealth inequality in their social network, this propels further inequality through reduced cooperation and reduced social connectivity. [SEE LETTER P.426](#)

SIMON GÄCHTER

Inequality is a growing concern in many societies¹. Like most important social phenomena, it is a complex issue that has many interacting sources and consequences^{1–3}. To understand inequality and its dynamics over time, multiple theoretical and empirical approaches are necessary. In this issue, Nishi *et al.*⁴ (page 426) use

laboratory-style experiments (conducted online) to study how the visibility of wealth inequality in people's social environment shapes the behavioural dynamics of inequality. The attraction of an experimental approach is that it allows the control of factors that are inherently uncontrollable in naturally occurring data. Crucially, for example, the experimenter can control the initial level of inequality and see how inequality evolves as a

function of people's behaviour alone^{5,6}.

Nishi and colleagues' experimental model used an assessment of people's willingness to contribute to public goods to test how initial wealth inequality and the structure of the social network influence the evolution of inequality. The researchers were particularly interested in the role of visibility of wealth — can mere observation of your neighbour's wealth lead to more inequality over time, even if such information does not change economic incentives? Visible wealth might have a psychological effect by triggering social comparisons and thereby influencing economic choices that have repercussions for inequality³.

In their online laboratory, the researchers endowed all participants with tokens, worth real money. The endowment differed across individuals and treatments: in a treatment without inequality, all participants initially received the same number of tokens; in a low-inequality treatment, participants had similar but different initial endowments; and in the high-inequality treatment there was a substantial starting difference between participants.

The groups typically comprised 17 people arranged at random in a social network in which, on average, about 5 people were linked ('neighbours'). In each of the 10 rounds of the following game, participants had to decide whether to behave pro-socially ('cooperate') by reducing their own wealth by 50 tokens per connected neighbour to benefit each of them by 100 tokens, or to behave pro-selfishly ('defect') by keeping their tokens for themselves. These decisions had consequences for accumulated wealth levels and inequality. At the end of each round, the subjects learnt whether their neighbours had cooperated or defected and 30% of participants were given the opportunity to change their neighbour, that is, to either sever an existing link or to create a new one.

A crucial manipulation in this experiment was wealth visibility. Under invisible conditions, the participants could observe only their own accumulated wealth. Under visibility, they could see the accumulated wealth of their connected neighbours but not the whole network. Thus, there were six conditions in total: three levels of initial wealth inequality in each of the two visibility conditions.

The results are complex but illuminating. The authors find that, under high initial wealth inequality, visibility of neighbours' accumulated wealth increases inequality over time relative to the invisibility condition, although absolute inequality decreases over time under both visibility conditions. The reason for the relative increase under visibility is that inequality drops only moderately, whereas under invisibility the reduction in inequality is substantial. By contrast, in the case of initial wealth equality, inequality increases — similarly in both visibility conditions. Under moderate initial



Figure 1 | Wealth on display. Nishi *et al.*⁴ use an experimental game to show that, when people can see the wealth of others whom they are linked with in a social network, inequality increases and the number of social connections decreases.

inequality, visibility leads to a small increase in inequality relative to invisibility.

Visibility of wealth also leads to lower social welfare, as measured by overall wealth (Fig. 1). By the end of the experiment, total accumulated wealth was substantially larger in the three conditions with invisible wealth than in the three conditions with visible wealth. The reason for this is that cooperativeness was lower under the condition of visible wealth compared to invisible wealth, and there were fewer links in the social network.

The most striking insight from these findings is the effect of wealth visibility on the dynamics of inequality: conspicuous inequality breeds more inequality. Although visibility of wealth does not change economic incentives in this experimental scenario, it invites social comparisons that, for various reasons^{3,7} worth exploring further, undermine cooperation and diminish social ties. This observation adds to existing^{8,9}, but sparse, evidence that public information about individual pay-offs leads to more competition, which in a public-goods setting triggers more 'free-riding' by individuals (defecting when others cooperate), to improve their own pay-offs.

Nishi and colleagues' findings raise several

intriguing methodological questions for future studies. For example, how much influence does the social network and its rewiring have on the main results of this experiment? Modelling interactions using a social network is certainly realistic, but is it crucial for the emergence of visibility effects in inequality? Another question concerns the result that visibility of wealth matters much less under initial equality of wealth. This is surprising, given that inequality of wealth increases over time and visibility effects should kick in, according to the results from the treatments with initial inequality. It is possible that these experiments, which used only ten iterations, might have been too short to allow for visibility effects arising as inequality grows.

The results also suggest substantive questions worthy of further research. As well as understanding the role of visibility of wealth (or pay-offs more generally) for cooperation, it would be interesting to gather evidence about how people's pro-social attitudes are affected by the ever-increasing amount of information about other people's consumption (as a signal of their wealth)¹⁰, which nowadays is spread on an almost global scale by social media. And how do visibility and social comparisons affect the dynamics of inequality when the relevant game is not one of cooperation but of competition? This is interesting because, in many interactions in our modern societies, not only initial endowments (wealth) matter but also resources that are allocated as people compete for scarce rewards — good jobs, for instance¹¹.

These are just some questions that can be investigated with the experimental model put forward by Nishi and colleagues. Their most general contribution is to showcase the power of experiments to contribute to our understanding of the behavioural dynamics of inequality. ■

Simon Gächter is in the Centre for Decision Research and Experimental Economics, University of Nottingham, Nottingham NG7 2RD, UK.
e-mail: simon.gaechter@nottingham.ac.uk

1. Atkinson, A. B. *Inequality: What Can Be Done?* (Harvard Univ. Press, 2015).
2. Chin, G. & Culotta, E. *Science* **344**, 818–821 (2014).
3. Frank, R. H. *Falling Behind: How Rising Inequality Harms the Middle Class* (Univ. California Press, 2013).
4. Nishi, A., Shirado, H., Rand, D. G. & Christakis, N. A. *Nature* **526**, 426–429 (2015).
5. Sadrieh, A. & Verbon, H. A. A. *Eur. Econ. Rev.* **50**, 1197–1222 (2006).
6. Gächter, S., Mengel, F., Tsakas, E. & Vostroknutov, A. <http://dx.doi.org/10.2139/ssrn.2351717> (2014).
7. Dohmen, T., Falk, A., Fliessbach, K., Sunde, U. & Weber, B. *J. Publ. Econ.* **95**, 279–285 (2011).
8. Huck, S., Normann, H.-T. & Oechssler, J. *Int. J. Industr. Organiz.* **18**, 39–57 (2000).
9. Nikiforakis, N. *Games Econ. Behav.* **68**, 689–702 (2010).
10. Veblen, T. *The Theory of the Leisure Class* (Macmillan, 1899).
11. Hopkins, E. & Kornienko, T. *Am. Econ. J.: Microeconomics* **2**(3), 106–137 (2010).



Cover illustration
Nik Spencer

Editor, *Nature*
Philip Campbell

Publishing
Richard Hughes

Insights Editor
Ursula Weiss

Production Editor
Elizabeth Batty

Art Editor
Nik Spencer

Sponsorship
Stephen Brown
Samantha Morley

Production
Ian Pope

Marketing
Hannah Phipps

Editorial Assistant
Rebecca White

The Macmillan Building
4 Crinan Street
London N1 9XW, UK
Tel: +44 (0) 20 7833 4000
e: nature@nature.com



nature publishing group

From wearable activity trackers to metagenomic sequencing and direct-to-consumer genetic testing, we are able to monitor our personal environment and health more than ever before. Rapid improvements in technology have also driven genetic discovery in human disease. Although the interpretation and definition of clinically relevant genetic variation remains a challenge, established examples are being used to stratify subgroups of patients to identify optimal treatments. Precision medicine is emerging as a natural extension that integrates research disciplines and clinical practice to build a knowledge base that can better guide individualized patient care.

At a crucial time when efforts such as the UK 100,000 Genomes Project and the US Precision Medicine Initiative seek to scale up population-based genome sequencing and integrate it with clinical data, we present this collection of reviews that assess progress towards the implementation of precision medicine. A new framework is required to bring together researchers, clinical laboratories, clinicians and patients in what Samuel Aronson and Heidi Rehm term a “precision-medicine ecosystem”.

Tailoring treatment to the patient is central: for this, Mary Relling and William Evans examine progress in pharmacogenomics and Luigi Naldini describes developments in gene therapy. Drug-development pipelines are benefiting from genomics in target-validation approaches that promise improved success rates and reduced costs, whereas clinical trials need to be redesigned to match the right trial to the right patient, as proposed by Andrew Biankin and colleagues.

For the next steps, the need to share data is greater than ever. The Global Alliance for Genomics and Health offers innovative solutions that facilitate research and clinical diagnoses while maintaining the privacy of sensitive data. Also essential will be a continually updated knowledge base to aid the interpretation of genetic tests. Although the challenges might seem formidable, we are encouraged by the early examples discussed here and are optimistic for the realization of precision medicine.

Nature is pleased to acknowledge the financial support of AstraZeneca in producing this Insight. As always, *Nature* carries sole responsibility for all editorial content.

Orli Bahcall
Senior Editor

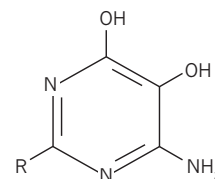
CONTENTS

REVIEWS

- 336 Building the foundation for genomics in precision medicine**
Samuel J. Aronson & Heidi L. Rehm

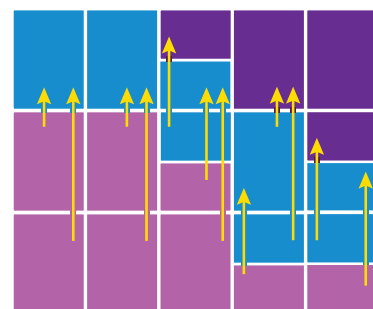


- 343 Pharmacogenomics in the clinic**
Mary V. Relling & William E. Evans



- 351 Gene therapy returns to centre stage**
Luigi Naldini

- 361 Patient-centric trials for therapeutic development in precision oncology**
Andrew V. Biankin, Steven Piantadosi & Simon J. Hollingsworth



Building the foundation for genomics in precision medicine

Samuel J. Aronson^{1,2} & Heidi L. Rehm^{1,3,4,5}

Precision medicine has the potential to profoundly improve the practice of medicine. However, the advances required will take time to implement. Genetics is already being used to direct clinical decision-making and its contribution is likely to increase. To accelerate these advances, fundamental changes are needed in the infrastructure and mechanisms for data collection, storage and sharing. This will create a continuously learning health-care system with seamless cycling between clinical care and research. Patients must be educated about the benefits of sharing data. The building blocks for such a system are already forming and they will accelerate the adoption of precision medicine.

The practice of medicine is an inexact science. The clinician assesses the patient's symptoms and decides which tests to perform to gather more data. They must determine the cause of the symptoms and the patient's prognosis, whether clinical intervention is warranted and, if so, which intervention to prescribe. To do this effectively, the clinician might need to assess several potential courses of action and incorporate all that is known about the patient into his or her decision. Human physiology is complex. In some cases, the cause of the patient's symptoms cannot be ascertained. In other cases, clinicians cannot gather enough data to make a fully informed decision. The guesswork inherent in the practice of medicine reduces the efficacy of the interventions that are prescribed.

Genetics is an important contributor to this complexity. Distinct genetic variants cause conditions that respond to different treatments yet share a similar set of symptoms. Without a mechanism to determine the underlying genetic cause of a set of symptoms, it might not be possible to determine which treatment will be most effective a priori. For instance, although there are many causes of lung cancer, only people who have an alteration in the gene *EGFR* respond to treatment with tyrosine kinase inhibitors^{1,2}. Similarly, many genetic lesions lead to a thickened heart and an increased risk of sudden cardiac death, but only people with mutations in the gene *GLA* respond to enzyme replacement therapy³. Even when the cause of a condition is known, unrelated genetic variants can affect treatment efficacy by altering the way in which drugs are metabolized or by increasing the likelihood of adverse events. For example, patients who are treated with conventional doses of the immunosuppressive drug azathioprine for an extended period are at risk of developing life-threatening myelosuppression if they harbour genetic variants that prevent the drug from being properly metabolized⁴. And approximately 6% of European populations carry *HLA-B* alleles that predispose them to potentially life-threatening hypersensitivity reactions if they are treated with the antiretroviral drug abacavir⁵.

Understanding the patient's genetic make-up is crucial for providing optimal care for many diseases. Clinicians now have access to an increasing array of tests that allow them to determine which genetic variants exist in their patients. These include: genotyping tests that look at variants in a patient's DNA sequence that are known to associate strongly with important clinical effects; panel-based gene sequencing, which looks at many genes related to a specific indication to detect known and new variants; sequencing of the exome — all known protein-coding

genes; and whole-genome analysis that attempts to sequence a patient's genome. However, simply determining which variants are present is insufficient. The implications of these variants must also be determined for each clinical indication. This genetic understanding must then be considered in conjunction with other clinical data to decide the path that will produce the best results for the patient.

The precision-medicine ecosystem

The goal of precision medicine is to enable clinicians to quickly, efficiently and accurately predict the most appropriate course of action for a patient. To achieve this, clinicians are given tools — in the form of tests and information-technology support — that are both compatible with their clinical workflow and economically feasible to deploy in the modern health-care environment. These tools help to simplify the process of managing the extreme biological complexity that underlies human disease. To support the creation and refinement of these tools, a precision-medicine 'ecosystem' is developing. This ecosystem is beginning to link clinicians, laboratories, research enterprises and clinical-information-system developers together in new ways. There is increasing hope that these efforts will create the foundation of a continuously learning health-care system that is capable of fundamentally accelerating the advance of precision-medicine techniques.

Interpretation is key to the precision-medicine ecosystem. It occurs at several levels. Individual variants can be interpreted in relation to specific indications. Sets of variants can be assessed in relation to their collective impact on patients. Genetic and clinical data can be combined to determine the best course of action for a patient. The quality of these interpretations is highly dependent on the data on which they are based. For this reason, research and clinical databases provide the foundation for precision medicine. Continuous learning in health care is in many ways driven by improvements to the content and structure of these resources. For example, the highly collaborative Clinical and Functional Translation of CFTR (CFTR2) project brought together researchers from around the world to share extensive patient data on *CFTR* variants to distinguish between pathogenic and benign lesions. This work is leading to more effective treatments for patients with cystic fibrosis⁶. Similarly, the Evidence-based Network for the Interpretation of Germline Mutant Alleles (ENIGMA) consortium has engaged a large collaborative network to share data on the *BRCA1* and *BRCA2* genes that predispose patients to breast and ovarian cancer⁷. Patients

¹Partners HealthCare Personalized Medicine, Boston, Massachusetts 02115, USA. ²Partners HealthCare Research Information Services and Computing, Charlestown, Massachusetts 02129, USA.

³Department of Pathology, Brigham & Women's Hospital, Boston, Massachusetts 02115, USA. ⁴Harvard Medical School, Boston, Massachusetts 02115, USA. ⁵The Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA.

with pathogenic variants in these genes can now take preventive action through monitoring and prophylactic surgical procedures⁸ and those with active breast cancer are candidates for targeted treatments^{9,10}. In addition, large collaborative programmes led by the US National Institutes of Health (NIH)-supported Clinical Genome (ClinGen) Resource¹¹ and the Global Alliance for Genomics and Health have begun to tackle the development of reliable resources for systematically defining the pathogenicity of all human variation through broad and targeted efforts.

When optimized, the infrastructure that supports the precision-medicine ecosystem efficiently manages and integrates the flow of material, knowledge and data needed to generate, validate, store, refine and apply clinical interpretations (Fig. 1). Biobanks link samples with patient data to support discovery. Research databases record the data, calculations and results that provide evidence for clinical interpretations. Clinical-knowledge-sharing networks enable the refinement of interpretations. Clinical laboratories and their information systems facilitate the consolidation of interpretations into reports and alerts. Electronic health records (EHRs) and associated systems help clinicians to apply results, both when they are received and as the patient's condition and knowledge of the variants evolve. Patient-facing infrastructure or 'portals' provide individuals with access to their genetic data and — if appropriate — the ability to decide how they should be used, including whether to participate in research. At present, much of this infrastructure is at a very early stage of development. However, the infrastructural foundation for precision medicine is beginning to emerge. In this Review, we explore its crucial components.

The patient viewpoint

The role of the patient in supporting precision medicine is becoming increasingly important. Patients are obtaining a growing number of genetic results in the course of their care. Typically, clinicians involved in their treatment order such tests for them. However, patients are also now able to access direct-to-consumer testing, sometimes through the help of someone who is not directly involved in their care. To ensure that precision medicine is tailored to the unique genetic make-up of each patient, we must gather as much information as possible from individual patients. Yet there are risks associated with widespread sharing of patient data. To gain access to these data, researchers must actively engage patients, teach them about the benefits of data sharing and help them to weigh up the risks and benefits. This can be done by making the process of obtaining consent more effective.

There are two major forms of consent that are relevant: consent for receiving medical treatment or procedures; and consent for releasing data or samples for use in research. In both cases, the risks and benefits must be conveyed to the patient. However, the conventional distinction is that obtaining consent for treatment focuses on benefits to the individual whereas obtaining consent for research focuses on generalizable knowledge¹². Increasingly, the line between clinical care and research is blurring; participation in research studies can lead to a direct improvement in outcome for the patient^{13,14}, and the continuous capture of clinical-care data has been proved an effective way to inform generalizable knowledge¹⁵. As a result, efforts are under way to ask all patients who enter the clinical-care setting to sign a form that permits their data to be used in research^{16–19}. In addition, those signing clinical genetic-testing consent forms now commonly agree to share their data broadly to help advance knowledge¹¹. Nevertheless, there is still a need for more uniform consenting processes. It is difficult to generate consent forms in language that is both easy to understand and robustly conveys the main issues associated with genetic testing. Sharing such language across institutions could be helpful in this context. Harmonizing consent language across providers, laboratories and biobanks would make it easier to administer and adhere to those agreements. Recently, the Regulatory and Ethics Working Group of the Global Alliance for Genomics and Health published a framework for the responsible sharing of genomic and health-related data²⁰. The group has also created

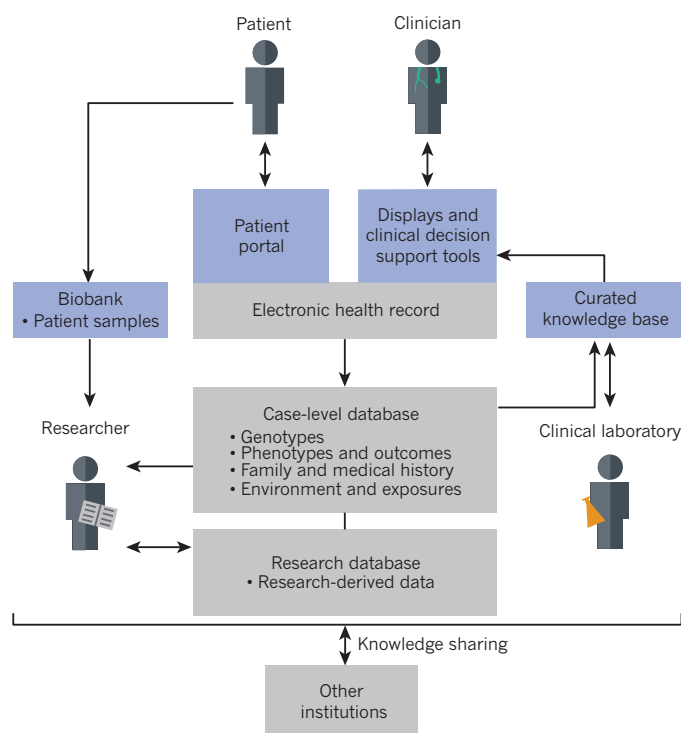


Figure 1 | The precision-medicine ecosystem. The precision-medicine ecosystem contains building blocks that optimally connect patients, clinicians, researchers and clinical laboratories to one another. Patients and clinicians access information through portals or EHRs. The ecosystem can include displays or CDS augmented by curated knowledge that is supplied and shared by multiple stakeholders. Case-level databases and biobanks receive case data and samples from clinical and research workflows. Researchers benefit from all of these information sources and also contribute to knowledge sources. Clinical laboratories leverage data and inform the clinical community as they assess genomic variation and its impact on human health.

consent tools and policies to aid the development of standardized approaches to obtaining consent and that support data sharing in the global community. Consistent with the Global Alliance for Genomics and Health framework, ClinGen has developed standardized consenting approaches (<http://clinicalgenome.org/data-sharing/>) for use in the clinical-care setting, which will enable sharing of genetic-test results and accompanying phenotypic data in the absence of research-study enrolment.

Some patients are extremely interested in supporting research and are willing to take proactive steps to facilitate the sharing of genetic information. The Global Network of Personal Genome Projects recruits volunteers who are prepared to share their genomic data and medical histories publicly. ClinGen manages the GenomeConnect patient portal, built on the Patient Crossroads platform, which allows individuals to share health and genetic information to form communities. The Platform for Engaging Everyone Responsibly (PEER), supported by the Genetic Alliance, enables individuals to control sharing, privacy and access preferences for their health and genomic data with a high degree of precision.

The clinician viewpoint

Clinicians gain access to patients' genetic information through tests. Tests have two components: a technical component that focuses on identifying which variants are present in the patient; and an interpretive component in which the implications of identified variants are assessed. In most scenarios, genetic testing is performed to determine either the cause of a specific indication or the most appropriate treatment²¹. However, exome and genome data can be reused to perform multiple

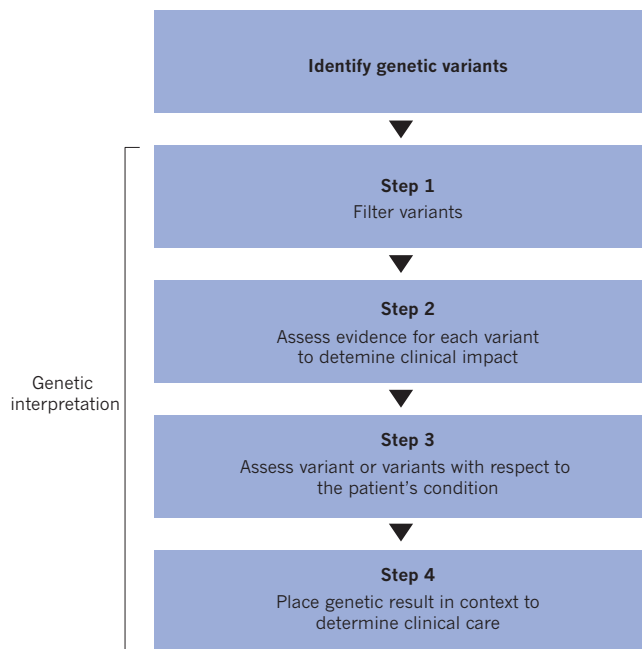


Figure 2 | Stages of the genetic interpretation process. Once genetic variants have been identified, they are filtered to select those of interest (step 1). Next, the evidence for each variant is assessed to determine the variant's clinical impact (step 2). One or more assessed variants are then interpreted with respect to the specific condition for which the patient is being investigated (step 3). Last, the overall genetic assessment is placed into the patient's clinical and personal context to inform the clinical-care decision-making process (step 4).

assessments over time. This opens up the possibility of obtaining and storing genome and exome sequences before disease manifests, with the intention that they will be interpreted and reinterpreted as indications arise. Irrespective of when the sequence is obtained, the interpretation step is crucial.

Clinical interpretation is a multiple-component process (Fig. 2). After a pool of variants have been identified, step 1 determines which of those variants should be assessed further. In step 2, the clinical impact of each of those variants is assessed. In step 3, the relevance of the combination of variants identified in step 2 are considered in relation to the patient's indication. Finally, in step 4, the test results are placed within the context of all known information about the patient to determine clinical care. Many laboratories base their reports on steps 1–3, although some simply report the variants they have identified. The patient's clinician often reviews the laboratory analysis and completes any remaining steps. These processes can be time consuming and therefore expensive. They also involve a considerable amount of professional judgement, which makes them subject to human error and differences of opinion. The quality and efficiency of these processes is highly dependent on the data available in clinical and research databases. For this reason, these databases are in many ways the core of the ecosystem that is needed to advance the practice of precision medicine. In addition, more standardized ways to evaluate evidence, such as those released by the American College of Medical Genetics and Genomics²², will be crucial for interpreting human genetic variation in a more consistent way. Finally, the open sharing of clinical interpretations to distribute the labour of variant assessment, to identify and resolve differences in interpretation, and to catalogue variation for research studies is essential for improving the care of patients with genetic-based conditions¹¹.

In each patient encounter, clinicians must address several questions that relate to precision medicine. First, they must assess whether genetics could be relevant — and if so, order the appropriate tests. Once the test results are received, the clinician must determine how to apply them. Then they must manage the results over time. Information-technology

support is needed to manage the large amount of patient data and other information that are required to execute these processes optimally. EHRs and their associated systems are the main means of providing such support to clinicians.

Electronic health records

EHRs are well positioned to be the apex of genetic information-technology support. They should serve as the clinician's gateway to all of the patient's information, including any genetic data. Information should be organized and displayed in a way that integrates with the clinician's workflow and facilitates diagnostic and treatment decisions. EHR and related systems can also provide clinicians with electronic clinical-decision support (CDS) that provides extra information about a genetic test or result through an e-resource or InfoButton^{23,24} that links to electronic resources such as websites or databases. They can also issue pre-test and post-test pharmacogenomic warnings that highlight potentially adverse interactions between drugs and specific genetic variants. Pre-test warnings are triggered when a clinician takes an action that should be informed by a genetic assessment but there is no record of the assessment being performed. Post-test alerts are triggered when an action is taken that may be contraindicated by a patient's genetic profile. An example is ordering a high dose of azathioprine for a patient with a thiopurine methyltransferase deficiency⁴. CDS systems can also alert clinicians when important information emerges on a patient's previously reported variant²⁵. In the future, CDS systems might be able to guide clinicians through complex scenarios that take into account multiple types of patient data, including genetics. Evolving such CDS is essential for the formation of a learning health-care system.

Genetic information and CDS do not necessarily have to be implemented directly into the EHR — it is possible to integrate EHRs with external systems^{25,26}. Such integration can be seamless so that clinicians need never know that they are working with multiple systems. Providing genetic support through the EHR is complex²⁷ and it is currently unclear how much genetic functionality EHR vendors will build into their systems. Some have indicated that they are unlikely to store full genomic sequencing in the EHR, instead choosing to link to external genomic data stores and focus their internal functionality on managing test results that have been interpreted at a higher level. Irrespective of whether patient genetic profiles are stored in the EHR, a constellation of systems will need to be tightly integrated with the EHR to provide optimal support to clinicians.

Displays of genetic information and CDS are often impossible to provide without robust access to the patient's genetic results and reports. The EHR and related clinician-facing systems must obtain genetic results from laboratory systems. This requires interfaces between the EHR provider and the laboratory. Creating these interfaces often involves establishing electronic connections that span multiple organizations and integrate systems from competing vendors. Relatively few such interfaces exist, largely because of the expense associated with creating them. Generally, results are transmitted from the laboratory to the provider by fax, which makes it difficult to keep the results organized in the EHR. CDS usually relies on access to structured electronic data that cannot be reliably extracted from a fax. Even if a genetic result is recorded in a structured format, this structure is often lost when the result is transferred to clinicians involved in the patient's care who operate out of different institutions. Any results from direct-to-consumer testing are also unlikely to be transferred in a structured format.

Several groups are working to promote interconnectivity that would enable CDS systems that incorporate genetic information. The Institute of Medicine Roundtable on Translating Genomic-Based Research for Health established the Displaying and Integrating Genetics Information Through the EHR Action Collaborative (DIGITizE AC). DIGITizE AC brings together clinicians, laboratories, vendors, standards organizations, government agencies and patient representatives to increase support for genetics in the EHR. The group has defined a set of genetics-based CDS rules that it seeks to roll out widely. This

will involve leveraging the frameworks of standards bodies, such as Health Level Seven (HL7) International, and the Logical Observation Identifiers Names and Codes (LOINC) database, as well as ontology and rule creators such as the Clinical Pharmacogenomics Implementation Consortium (CPIC). The National Human Genome Research Institute (NHGRI) has also established several consortia with EHR working groups to investigate how genetics can be supported in the EHR. These include: the Electronic Medical Records and Genomics (eMERGE) Network²⁸, Clinical Sequencing Exploratory Research (CSER) and ClinGen¹¹.

However, the problem that these organizations are trying to solve is very difficult. Standardized message formats and ontologies are the best way to reduce the cost of establishing the laboratory–provider and provider–provider interfaces needed to underlie precision medicine. However, these standards are helpful only if they robustly account for the different real-world scenarios they are intended to support and are broadly implemented by the vendor community (Fig. 3). Developing such standards requires an enormous amount of input from groups that combine deep clinical, laboratory, vendor and information-technology expertise. The DIGITize AC has found that even defining the specific requirements for its initial set of narrow-use cases entailed a considerable amount of interdisciplinary effort. Much more work is needed to build truly robust, general-purpose standards.

The clinical laboratory viewpoint

Clinical laboratories sit at the core of the interpretative process. Ideally, they provide both the evidence for individual variants as well as a case-level report that places all potentially relevant variants in the context of the patient's presentation. Laboratories that perform genome sequencing often discover variants that they have never seen before, which must then be assessed. Similarly, variants that have been seen before might need to be reassessed as new knowledge emerges. Variant assessment is becoming an important factor in the cost of genetic tests. It must be performed by skilled individuals because errors could result in inappropriate patient care. Yet we know that variants can be interpreted differently. As of 11 September 2015, 369 organizations had submitted a total of 158,668 variants to ClinVar, a National Center for Biotechnology Information (NCBI) database that acts as a single centralized public repository to which institutions can submit their interpreted variants as well as retrieve data from others²⁹. At least 2,000 of these have been interpreted differently by submitters¹¹.

Laboratories and clinicians can be assisted in two ways: better access to variant assessments performed by other institutions using consistent approaches, and tools to improve and standardize the variant assessment process.

Building clinical genomic knowledge

Sharing variant- and gene-level assessments between laboratories and clinicians can increase the quality and efficiency of the variant assessment process. Multiple efforts are under way to increase the sharing of such knowledge^{30–34}. The ClinGen programme is building an authoritative central resource that defines the clinical relevance of genomic variants for use in precision medicine and research. The programme aims are to increase the rate of submission to ClinVar and to improve the content of ClinVar and other genomic resources through expert curation. ClinGen has worked together with ClinVar to create a 'star system' that defines the level of review for each variant that is submitted to ClinVar¹¹. ClinGen working groups have been established in multiple clinical domains to curate gene–disease relationships and to interpret variants through expert consensus.

Centralized knowledge repositories can also be created by linking together the infrastructure that supports different laboratories. For example, laboratories that use the GeneInsight Lab application³⁵ are able to use the system to communicate and share knowledge in real time. This functionality has been used to create a network called VariantWire and also supports the Canadian Open Genetics Repository

(COGR)³⁶ network of Canadian labs. Importantly, an organization can both participate in a knowledge-sharing network and contribute their data to ClinVar. By adopting a standardized infrastructure that helps to structure data for submission to ClinVar, public sharing becomes cheaper, more efficient and more comprehensive with respect to supplying the supporting evidence.

Case repositories and biobanks

An important driver of improvements to variant assessment processes is the collection and analysis of case data. Clinical and research laboratories often develop case repositories. The power of these repositories is a function of the number of cases that they contain. Therefore sharing cases across institutions is beneficial. However, it is difficult to combine data that have been stored in information systems developed by different groups. Trade-offs must be made when deciding what data to capture and how deeply to standardize and structure them. The amount of data in a case repository can be increased by allowing contributors to deposit heterogeneous data that are incomplete or inconsistently validated and may therefore be difficult to process downstream³⁷. If repository developers insist on the submission of complete, validated and consistent data, many cases will have to be excluded.

Several databases have been launched that share case-level data across broad disease areas. The NCBI's database of Genotypes and Phenotypes (dbGaP)³⁸ places minimal restrictions on the types of case data that can be submitted and therefore serves as a generalized repository. However, because phenotypic data are often limited, making informative use of the information is difficult. Similarly, the European Bioinformatics Institute (EBI) maintains the European Genome-phenome Archive for storing case-level genomic data. The International Cancer Genome Consortium (ICGC)³⁹ and The Cancer Genome Atlas (TCGA)⁴⁰ have each set up large repositories of somatic cancer sequencing data. The American Society of Clinical Oncology (ASCO) is looking to incorporate the tracking of patient outcomes to enable a learning health-care system in its CancerLinQ platform⁴¹. Repositories have also been developed through direct patient participation and span non-profit, academic and commercial activities.

Access to clinical specimens associated with patient data is often

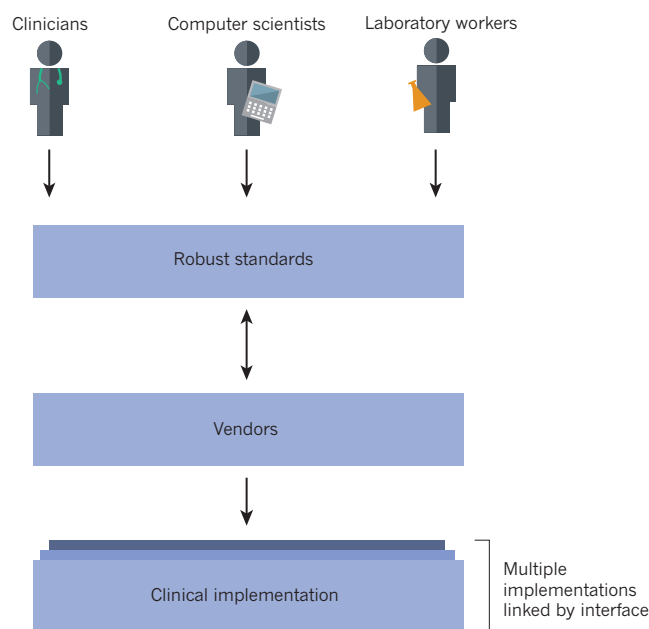


Figure 3 | Creating and implementing robust standards for the description and structuring of data in laboratory processing and patient-care systems. Professionals with diverse expertise interact with vendors of laboratory-information systems and EHR systems to iteratively design and implement standards that effectively enable techniques to be used in the clinic.

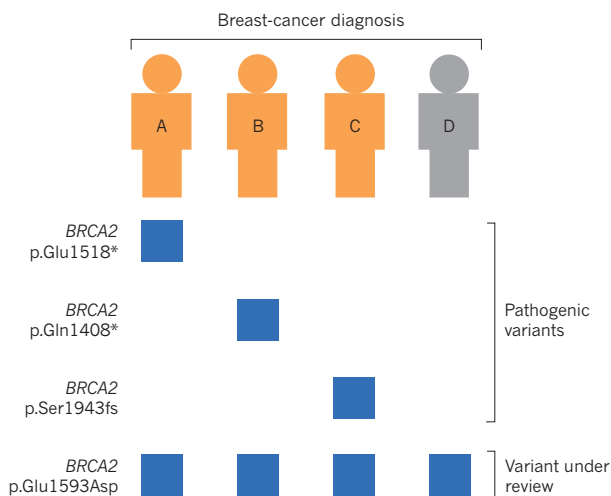


Figure 4 | Example of a learning health-care system. Case data can be shared between laboratories to support variant assessment. In this example, the *BRCA2* p.Glu1593Asp variant in case D is classified initially as being of ‘uncertain significance’. After accessing genetic and phenotypic patient data from cases A, B and C, in which there are other genetic explanations for the clinical phenotype, the necessary evidence becomes available to classify the *BRCA2* p.Glu1593Asp variant as ‘likely benign’.

necessary to fully inform discovery and continuous learning. The costs associated with collecting samples prospectively for research studies are enormous. However, when samples are collected in bulk and then placed into biobanks, which allows their reuse across studies, the costs decrease considerably. One of the keys to success is ensuring that strong consent processes are in place. Biobanks are moving from repositories of de-identified, unconsented specimens towards fully consented models that allow association with longitudinal health records. This is another area in which universal consent language would simplify the infrastructure development process. Direct engagement of patients in this type of sharing can help participants to balance the risks and benefits. Patients must fully understand the scope of their consent and be actively informed and engaged as participants in the advancement of knowledge. The return to patients of clinically relevant results from research studies should also be supported because this allows patients to benefit directly from their data sharing.

In addition to supporting individual case-level repositories and biobanks, efforts are being made to connect divergent databases. Launched in 2015, the Matchmaker Exchange is a centralized network for sharing case-level data within an international set of case-level repositories focused on gene discovery⁴². Although each database has its own data schema, the development of a common application programming interface⁴³ means that users can query genomic and phenotypic data across multiple systems. This has encouraged the member databases to move towards implementing a common set of fields to facilitate effective data exchange for gene discovery.

There is no doubt that other such efforts will emerge, particularly as the Genomics England and Precision Medicine Initiative programmes develop. Furthermore, the data-standardization efforts that help to establish interfaces between laboratories and providers could assist the development of these case repositories. Organizations that supply patient data to these efforts must develop mechanisms for collecting data more uniformly and for sharing them consistently. By enabling patients to contribute their data directly, the collection of phenotype data can be accelerated and broadened⁴⁴.

Knowledge resources and tools

In addition to clinical-knowledge and case-sharing networks, many laboratories and clinicians use research-grade knowledge resources and tools. Many types of tools and resources are used daily in the clinical

workflow, even if they are not intended for direct clinical usage. *In silico* assessment tools use computational algorithms to assess the likely effects of DNA variation^{45–47}. Genome browsers can display multiple tracks of information, including species conservation data, the location of gene transcripts and regulatory elements, and population genetic variation. A number of allele-frequency databases, such as dbSNP Short Genetic Variations, the National Heart, Lung, and Blood Institute (NHLBI) Grand Opportunity Exome Sequencing Project’s Exome Variant Server, the Exome Aggregation Consortium (ExAC) Browser and the 1000 Genomes Project⁴⁸, provide data that are used by clinical laboratories to define variation that is unlikely to cause Mendelian disorders. The NCBI’s PubMed database provides access to published biomedical literature, and public and commercialized efforts exist to curate and present the data contained in such literature in a more useful format, such as Online Mendelian Inheritance in Man (OMIM) and The Human Gene Mutation Database⁴⁹. To define the role of these knowledge resources, users must assess which resources are useful, how their quality is controlled and how best they can be integrated into clinical workflows⁵⁰. ClinGen maintains a list (<http://clinicalgenome.org/tools/web-resources/>) of web-based tools that members of its community have found useful. This resource has also been designed to serve as an e-resource that can be accessed through EHRs.

The researcher viewpoint

The advent of precision medicine and its supporting infrastructure has given researchers the ability to influence clinical care directly. The release of innovative research tools and the addition of new information in the form of knowledge bases can immediately influence patient care by changing how genomic variants are assessed. Although data obtained through clinical settings typically require processing in specific ways, those obtained from research tend to be more flexible. This means that researchers can often tolerate more variability and inconsistency in their data sources than clinicians. For this reason, new infrastructure is usually released to the research community before it is optimized for clinical use.

Clinical knowledge-sharing infrastructure and case repositories, especially when combined with EHR-derived content, can provide clinicians and clinical laboratories not only with unprecedented access to clinical data, but also make this information accessible to researchers. New models are emerging for the broad sharing of data for discovery purposes. For example, the crowd-sourced approach to solving complex biological problems taken by the DREAM Challenges is now being applied to clinical-trial data in the hope of advancing precision medicine^{51,52}. Such challenges pair the brightest computer scientists with unprecedented open access to data to allow the development of highly informative models for predicting patient outcomes. To support this era of open data and discovery, it is crucial that appropriate consent approaches are established to allow clinical data to be used in these ways. Opportunities for discovery are also being created as the cost of sequencing falls and the processing capabilities of ‘big data’ become increasingly accessible⁵³. All these factors could contribute further to the formation of a continuously learning health-care system that simultaneously engages clinical-care providers and researchers and is necessary to support the development of precision medicine.

Realizing continuously learning health care

Ideally, continuous learning in health care would involve the capture of all incremental data, knowledge and experience gained through each patient interaction. This information would then be used in real time to improve the care of current and future patients. The ability to stratify patients, understand scenarios and optimize decision-making would consistently improve based on the myriad data obtained during the care-delivery process. This would be the ultimate expression of precision medicine. The infrastructure we discuss in this Review represents initial steps in this direction. We have already seen evidence to show that continuous-learning processes are achievable. Figure 4

depicts an example of a continuous-learning system in which hypothetical historical patient data from breast-cancer testing are accessed to determine the pathogenicity of a new *BRCA1* variant as 'likely benign'. The variant would otherwise be considered of 'uncertain significance'. Clinical laboratories often classify variants on the basis of historical clinical case histories. Because each laboratory has access to only a fraction of patients tested, optimal learning can only happen when data are shared broadly between organizations. ClinGen has made advances in clinical-laboratory genetic-data sharing through the use of the ClinVar database¹¹. However, this level of sharing is more likely to lead to therapeutic development and improved outcomes if the results of genetic testing are accompanied by greater amounts of patient health data. This will require the emerging genetic infrastructure to be extended, such that it can integrate with as many other forms of patient data as possible.

Addressing barriers to precision medicine

Multiple issues must be overcome for personalized medicine to reach its potential, as summarized by Joyner and Paneth in seven key questions⁵⁴. Although some doubt has been expressed that personalized medicine will reach its full potential for common diseases, the recent shift in emphasis to studies of the genetic basis of rare diseases and somatic cancer could provide tangible success in this field. For example, mechanistic understanding of rare disease and cancer pathways might inform the understanding of common diseases and approaches to reducing risk more effectively than has been achieved through genome-wide association studies. However, to ensure that we can learn from our evolving experience in the diagnosis and treatment of all types of disease, continuously learning health-care systems and broad data-sharing approaches must be supported. The absence of such systems is likely to be responsible for the limited success of personalized medicine to date. The continuous-learning infrastructure could be used to add a testing methodology for new hypotheses, in which real-time evaluation is repeatedly conducted against a limited set of treatment decisions for a given condition to determine which treatments provide the best results for different patient subgroups. Improved decision-making — at present, based on access to more up-to-date knowledge, and in the future, based on real-time evaluation techniques — has the potential to partially offset cost concerns by reducing expenditures associated with unnecessary or ineffective care. These improvements are also likely to generate public-health benefits. Improved infrastructure to capture both test results and patient outcomes should enable the measurement of such benefits.

The type and quality of patient data stored in EHRs are clearly issues that need to be addressed to support a continuously learning health-care system. In our experience, the investment required to capture higher-quality and more clinically relevant data is made only when a near-term financial return on those investments can be established. An improved foundational infrastructure provides an expanded basis for innovation and thereby facilitates the development of tools and analysis that are capable of justifying these foundational investments.

Open data-sharing resources, as well as the principle of open data itself, can help to reduce the cost of conducting genetic research. They can also limit the number of conflict-of-interest problems that occur as academic medical centres increasingly partner with commercial activities. Although this infrastructure does not help to solve generalized funding issues, it does set a precedent for sharing data rather than keeping it proprietary. In doing so, it reduces the scope and impact of potential conflicts and helps to ensure that commercial relationships are based on open principles.

Future directions

Despite compelling examples of the use of genomics to support precision medicine, the core building blocks that will be necessary to scale up the field are still in a very primitive state. However, as the community works to improve these building blocks and link them up, transformations are beginning to occur. Clinicians, researchers, laboratories and

vendors are working together to build the tools that will close the distance between each stakeholder. It is becoming easier to move, compare, apply and reproduce knowledge, data and samples. The basic infrastructure required to support a continuously learning health-care system has started to evolve spontaneously in many different areas. Furthermore, a cultural change is emerging as researchers, clinicians and patients embrace the open sharing of data to facilitate scientific advancement. Although it is unclear how long it will take to build an infrastructure that fully supports the widespread sharing and effective use of genomic and health data, the ultimate result will be a transformation of health care that allows continuous advances in medicine to occur within a clinical-care system that is less dependent on externally funded research endeavours. ■

Received 6 June; accepted 11 August 2015.

- Lynch, T. J. *et al.* Activating mutations in the epidermal growth factor receptor underlying responsiveness of non-small-cell lung cancer to gefitinib. *N. Engl. J. Med.* **350**, 2129–2139 (2004).
 - Paez, J. G. *et al.* *EGFR* mutations in lung cancer: correlation with clinical response to gefitinib therapy. *Science* **304**, 1497–1500 (2004).
 - Morel, C. F. & Clarke, J. T. The use of agalsidase alfa enzyme replacement therapy in the treatment of Fabry disease. *Expert Opin. Biol. Ther.* **9**, 631–639 (2009).
 - Relling, M. V. *et al.* Clinical Pharmacogenetics Implementation Consortium guidelines for thiopurine methyltransferase genotype and thiopurine dosing: 2013 update. *Clin. Pharmacol. Ther.* **93**, 324–325 (2013).
 - Martin, M. A. *et al.* Clinical pharmacogenetics implementation consortium guidelines for HLA-B genotype and abacavir dosing: 2014 update. *Clin. Pharmacol. Ther.* **95**, 499–500 (2014).
 - Cutting, G. R. Cystic fibrosis genetics: from molecular understanding to clinical application. *Nature Rev. Genet.* **16**, 45–56 (2015).
 - Spurdle, A. B. *et al.* ENIGMA—evidence-based network for the interpretation of germline mutant alleles: an international initiative to evaluate risk and clinical significance associated with sequence variation in *BRCA1* and *BRCA2* genes. *Hum. Mutat.* **33**, 2–7 (2012).
 - Domchek, S. M. *et al.* Association of risk-reducing surgery in *BRCA1* or *BRCA2* mutation carriers with cancer risk and mortality. *J. Am. Med. Assoc.* **304**, 967–975 (2010).
 - Audeh, M. W. *et al.* Oral poly(ADP-ribose) polymerase inhibitor olaparib in patients with *BRCA1* or *BRCA2* mutations and recurrent ovarian cancer: a proof-of-concept trial. *Lancet* **376**, 245–251 (2010).
 - Tutt, A. *et al.* Oral poly(ADP-ribose) polymerase inhibitor olaparib in patients with *BRCA1* or *BRCA2* mutations and advanced breast cancer: a proof-of-concept trial. *Lancet* **376**, 235–244 (2010).
 - Rehm, H. L. *et al.* ClinGen — The Clinical Genome Resource. *N. Engl. J. Med.* **372**, 2235–2242 (2015).
- This article describes ClinGen, an NIH-supported programme to build an authoritative central resource that defines the clinical relevance of genomic variants for use in precision medicine and research, employing systematic sharing of clinical knowledge and expert curation.**
- US Department of Veterans Affairs Office of Research & Development. *Informed Consent for Human Subjects Research: a Primer* http://www.research.va.gov/resources/pubs/docs/consent_primer_final.pdf (VA Boston Health Care System, 2002).
 - Jameson, E., Jones, S. & Wraith, J. E. Enzyme replacement therapy with laronidase (Aldurazyme®) for treating mucopolysaccharidosis type I. *Cochrane Database Syst. Rev.* **11**, CD009354 (2013).
 - Hacein-Bey Abina, S. *et al.* Outcomes following gene therapy in patients with severe Wiskott–Aldrich syndrome. *J. Am. Med. Assoc.* **313**, 1550–1563 (2015).
 - Murphy, S. N. *et al.* High throughput tools to access images from clinical archives for research. *J. Digit. Imaging* **28**, 194–204 (2015).
 - McCarty, C. A. *et al.* The eMERGE Network: a consortium of biorepositories linked to electronic medical records data for conducting genomic studies. *BMC Med. Genomics* **4**, 13 (2011).
 - Allen, N. L. *et al.* Biobank participants' preferences for disclosure of genetic research results: perspectives from the OurGenes, OurHealth, OurCommunity project. *Mayo Clin. Proc.* **89**, 738–746 (2014).
 - Toledo, J. B. *et al.* A platform for discovery: The University of Pennsylvania Integrated Neurodegenerative Disease Biobank. *Alzheimers Dement.* **10**, 477–484 (2014).
 - Milani, L., Leitsalu, L. & Metspalu, A. An epidemiological perspective of personalized medicine: the Estonian experience. *J. Intern. Med.* **277**, 188–200 (2015).
 - Knoppers, B. M. Framework for responsible sharing of genomic and health-related data. *HUGO J.* **8**, 3 (2014).
 - Korf, B. R. & Rehm, H. L. New approaches to molecular diagnosis. *J. Am. Med. Assoc.* **309**, 1511–1521 (2013).
 - Richards, S. *et al.* Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* **17**, 405–423 (2015).

These guidelines provide a standardized approach to the interpretation of genetic variants for monogenic disease.

23. Hoffman, M. A. & Williams, M. S. Electronic medical records and personalized medicine. *Hum. Genet.* **130**, 33–39 (2011).
 24. Del Fiol, G. *et al.* Integrating genetic information resources with an EHR. *AMIA Annu. Symp. Proc.* **2006**, 904 (2006).
 25. Aronson, S. J. *et al.* Communicating new knowledge on previously reported genetic variants. *Genet. Med.* **14**, 713–719 (2012).
 26. Starren, J., Williams, M. S. & Bottinger, E. P. Crossing the omic chasm: a time for omic ancillary systems. *J. Am. Med. Assoc.* **309**, 1237–1238 (2013).
 27. Kho, A. N. *et al.* Practical challenges in integrating genomic data into the electronic health record. *Genet. Med.* **15**, 772–778 (2013).
- This review summarizes challenges that the eMERGE consortium has encountered when integrating genetics into the EHR and suggests approaches for addressing these challenges.**
28. Gottesman, O. *et al.* The Electronic Medical Records and Genomics (eMERGE) Network: past, present, and future. *Genet. Med.* **15**, 761–771 (2013).
 29. Landrum, M. J. *et al.* ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* **42**, D980–D985 (2014).
 30. Bérout, C., Collod-Bérout, G., Boileau, C., Soussi, T. & Junien, C. UMD (Universal Mutation Database): a generic software to build and analyze locus-specific databases. *Hum. Mutat.* **15**, 86–94 (2000).
 31. Sosnay, P. R. *et al.* Defining the disease liability of variants in the cystic fibrosis transmembrane conductance regulator gene. *Nature Genet.* **45**, 1160–1167 (2013).
 32. Firth, H. V. *et al.* DECIPHER: Database of Chromosomal Imbalance and Phenotype in Humans Using Ensembl Resources. *Am. J. Hum. Genet.* **84**, 524–533 (2009).
 33. Miller, D. T. *et al.* Consensus statement: chromosomal microarray is a first-tier clinical diagnostic test for individuals with developmental disabilities or congenital anomalies. *Am. J. Hum. Genet.* **86**, 749–764 (2010).
 34. Thompson, B. A. *et al.* Application of a 5-tiered scheme for standardized classification of 2,360 unique mismatch repair gene variants in the InSiGHT locus-specific database. *Nature Genet.* **46**, 107–115 (2014).
 35. Aronson, S. J. *et al.* The GeneInsight Suite: a platform to support laboratory and provider use of DNA-based genetic testing. *Hum. Mutat.* **32**, 532–536 (2011).
 36. Lerner-Ellis, J., Wang, M., White, S. & Lebo, M. S. & Canadian Open Genetics Repository Group. Canadian Open Genetics Repository (COGR): a unified clinical genomics database as a community resource for standardising and sharing genetic interpretations. *J. Med. Genet.* **52**, 438–445 (2015).
 37. Riggs, E. R., Jackson, L., Miller, D. T. & Van Vooren, S. Phenotypic information in genomic variant databases enhances clinical care and research: the International Standards for Cytogenomic Arrays Consortium experience. *Hum. Mutat.* **33**, 787–796 (2012).
 38. Tryka, K. A. *et al.* NCBI's Database of Genotypes and Phenotypes: dbGaP. *Nucleic Acids Res.* **42**, D975–D979 (2014).
 39. Zhang, J. *et al.* International Cancer Genome Consortium Data Portal—a one-stop shop for cancer genomics data. *Database (Oxford)* **2011**, bar026 (2011).
 40. The Cancer Genome Atlas Research Network *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nature Genet.* **45**, 1113–1120 (2013).
 41. Schilsky, R. L., Michels, D. L., Kearbey, A. H., Yu, P. P. & Hudis, C. A. Building a rapid learning health care system for oncology: the regulatory framework of

CancerLinQ. *J. Clin. Oncol.* **32**, 2373–2379 (2014).

This article provides an overview of the challenges of applying precision medicine techniques to cancer and then describes the CancerLinQ system and the regulatory framework under which it operates.

42. Philippakis, A. A. *et al.* The matchmaker exchange: a platform for rare disease gene discovery. *Hum. Mutat.* <http://dx.doi.org/10.1002/humu.22858> (2015).
- This paper describes an international system for sharing genomic cases to aid in gene discovery.**
43. Buske, O. J. *et al.* The matchmaker exchange API: automating patient matching through the exchange of structured phenotypic and genotypic profiles. *Hum. Mutat.* <http://dx.doi.org/10.1002/humu.22850> (2015).
 44. Almalki, M., Gray, K. & Sanchez, F. M. The use of self-quantification systems for personal health information: big data management activities and prospects. *Health Inf. Sci. Syst.* **3** (suppl.), S1 (2015).
 45. Thusberg, J., Olatubosun, A. & Vihinen, M. Performance of mutation pathogenicity prediction methods on missense variants. *Hum. Mutat.* **32**, 358–368 (2011).
 46. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nature Genet.* **46**, 310–315 (2014).
 47. Jian, X., Boerwinkle, E. & Liu, X. In silico tools for splicing defect prediction: a survey from the viewpoint of end users. *Genet. Med.* **16**, 497–503 (2014).
 48. The 1000 Genomes Project Consortium. A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061–1073 (2010); erratum **473**, 544 (2011).
 49. Stenson, P. D. *et al.* The Human Gene Mutation Database: building a comprehensive mutation repository for clinical and molecular genetics, diagnostic testing and personalized genomic medicine. *Hum. Genet.* **133**, 1–9 (2014).
 50. Gargis, A. S. *et al.* Assuring the quality of next-generation sequencing in clinical laboratory practice. *Nature Biotechnol.* **30**, 1033–1036 (2012).
 51. Jarchum, I. & Jones, S. DREAMing of benchmarks. *Nature Biotechnol.* **33**, 49–50 (2015).
 52. Abdallah, K., Hugh-Jones, C., Norman, T., Friend, S. & Stolovitzky, G. The Prostate Cancer DREAM Challenge: A community-wide effort to use open clinical trial data for the quantitative prediction of outcomes in metastatic prostate cancer. *Oncologist* **20**, 459–460 (2015).
 53. O'Driscoll, A., Daugelaite, J. & Sleator, R. D. 'Big data', Hadoop and cloud computing in genomics. *J. Biomed. Inform.* **46**, 774–781 (2013).
- This review discusses cloud computing and big data concepts and their application to the field of genomics.**
54. Joyner, M. J. & Paneth, N. Seven questions for personalized medicine. *J. Am. Med. Assoc.* <http://dx.doi.org/10.1001/jama.2015.7725> (2015).

Acknowledgements H.L.R. was supported in part by NIH grants U41HG006834, U01HG006500 and U19HD077671. S.J.A. was supported in part by U41HG006834.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare competing financial interests: see go.nature.com/7jayhk for details. Readers are welcome to comment on the online version of this paper at go.nature.com/7jayhk. Correspondence should be addressed to H.L.R. (hrehm@partners.org).

Pharmacogenomics in the clinic

Mary V. Relling¹ & William E. Evans¹

After decades of discovery, inherited variations have been identified in approximately 20 genes that affect about 80 medications and are actionable in the clinic. And some somatically acquired genetic variants direct the choice of ‘targeted’ anticancer drugs for individual patients. Current efforts that focus on the processes required to appropriately act on pharmacogenomic variability in the clinic are moving away from discovery and towards implementation of an evidenced-based strategy for improving the use of medications, thereby providing a cornerstone for precision medicine.

Pharmacogenomics focuses on the identification of genetic variants that influence drug effects, typically through alterations in pharmacokinetics — that is, how the drug is absorbed, distributed, metabolized or eliminated — or pharmacodynamics, by modifying its target or by perturbing the biological pathways that shape a patient’s sensitivity to its pharmacological effects. In conditions apart from cancer and infectious diseases, the genetic variations of interest are primarily in germline DNA, and are inherited from a parent or are *de novo* changes that alter the function of gene products. In cancer, both inherited and somatically acquired variants can influence a patient’s response to treatments. In infectious diseases, genetic variation can affect a pathogen’s sensitivity to antimicrobial drugs¹. Advances in genome interrogation technology and in analytical approaches have facilitated the evolution of a discovery model from candidate gene studies towards agnostic genome-wide analyses of patient populations with specific drug-response phenotypes — for example, toxicity or desired pharmacological effects. In fact, current technologies for genome interrogation are sufficiently robust that defining the drug-response phenotype has become the more difficult component of pharmacogenomics research. Once a pharmacogenomic relationship has been discovered and validated, there are many obstacles to translating it into clinical practice. Such translation requires that effective, alternative therapy is available for those with ‘high-risk’ genotypes, as well as improvements to health-care systems, structured approaches to guide prescribing (for example, algorithms), and implementation of point-of-care electronic clinical decision support (CDS), to make it feasible to utilize genetics appropriately to guide drug prescribing.

A decade ago, we laid out a vision of how evolving genome technologies could be deployed to facilitate pharmacogenomic discoveries², and here in this Review, we extend this vision to address how these relationships can be translated into tools to optimize the use of medications in the clinic.

Evolution of pharmacogenomics

The earliest origins of pharmacogenomics are unclear; in 510 BC the Greek philosopher and mathematician Pythagoras reported that a subset of people who ingested broad beans (*Vicia faba*) experienced potentially fatal haemolytic anaemia (Fig. 1). Over two thousand years later, this reaction was attributed to an inherited deficiency in the enzyme glucose-6-phosphate dehydrogenase (G6PD), which also predisposes patients to haemolysis from several medications, including rasburicase and the antimalarial primaquine³. In 1909, while studying the common

bean *Phaseolus vulgaris*, the Danish pharmacist Wilhelm Johannsen coined the terms genotype and phenotype, linking genotype to the effects of volatile organics, a presage to pharmacogenetics⁴. A clustering of drug-metabolizing enzyme activities by racial group strongly suggested a genetic component to population variation^{5,6}.

In 1959, the German geneticist Friedrich Vogel was the first to use the term ‘pharmacogenetics’⁷, a concept that was bolstered by landmark studies by the pharmacologists Elliott Vesell and John G. Page showing that the pharmacokinetic profile of the pain-relieving drug antipyrine is more similar in monozygotic twins than in dizygotic twins⁸. The clinical relevance of the field was reinforced when family studies in different racial groups indicated that differences in the metabolism of isoniazid — a treatment for tuberculosis — and its side effect, peripheral neuritis, were inherited as an autosomal recessive trait^{9,10}. Decades later, differences in the metabolism of isoniazid were shown to be caused by inherited variants in the *NAT2* gene, which encodes the *N*-acetyltransferase 2 enzyme^{11,12}.

Other family studies conducted between the 1960s and 1980s documented patterns of inheritance for many drug effects, which eventually led to molecular studies that revealed the inherited determinants for many of the traits. In 1987, *CYP2D6* became the first polymorphic human drug-metabolizing gene to be cloned and characterized¹³. In the 1990s, the potential clinical utility of pharmacogenomics was clearly illustrated for several genes^{14,15}, including *TPMT*, which encodes the enzyme thiopurine methyltransferase. People with an inherited deficiency in this enzyme were found to experience haematopoietic toxicity on administration of the antileukaemic and immunosuppressive thiopurine drugs mercaptopurine and azathioprine¹⁶, although implementation of this finding in the clinic progressed slowly at that time¹⁷.

As in most areas of genetics, the rate of pharmacogenetic discovery has been accelerated by the Human Genome Project and by improved technologies for the genome-wide interrogation of variation. This shortened the timeline for discovery and enabled agnostic genome-wide studies of populations of patients with specific drug-effect phenotypes, often leading to the identification of unanticipated genetic variants that were statistically associated with drug effects. These genome-wide strategies also helped to bring the term ‘pharmacogenomics’ into the pharmacology lexicon¹⁸.

Discoveries that emerge from genome-wide or candidate-gene strategies require independent validation before they can be translated into clinical diagnostics. The validation process can be facilitated by elucidating the mechanisms that determine how the variation alters drug

¹Department of Pharmaceutical Sciences, St. Jude Children’s Research Hospital, Memphis, Tennessee 38105-2794, USA.

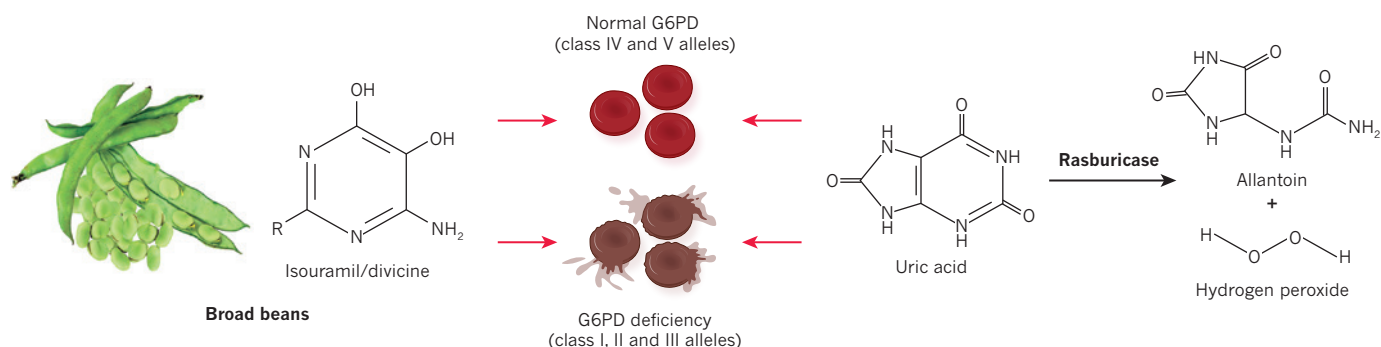


Figure 1 | Inherited G6PD deficiency and haemolysis. People with a deficiency in the enzyme glucose-6-dehydrogenase (G6PD) can develop haemolytic anaemia after eating broad beans (left) or when taking the drug rasburicase (right). In the case of broad beans, the reaction is the result of the chemical moieties isouramil ($R = OH$) and divicine ($R = NH_2$), and for

rasburicase it is the hydrogen-peroxide product of the breakdown of uric acid. Red blood cells of G6PD-deficient individuals — primarily those with the class II ‘Mediterranean’ allele of *G6PD* — produce insufficient NADPH to protect from oxidative damage caused by these moieties, which leads to haemolysis. Consequently, rasburicase is contraindicated in G6PD-deficient individuals¹¹².

responses. Genetic variants often differ according to ancestry, which can confound the translation of pharmacogenetic traits from one population to another. For instance, genetic polymorphisms in the genes *CYP2C9* and *VKORC1* have a population-specific influence on the anticoagulant effects of warfarin¹⁹. Furthermore, it is becoming increasingly evident that many drug effects are influenced by multiple variants in the same gene — some of which are rare — and by variants in multiple genes within the same patient. The Genomics England 100,000 Genomes Project and the US National Institutes of Health (NIH) Pharmacogenomics Research Network are two of many ongoing efforts to facilitate genome discoveries and their translation into diagnostics. Their findings may eventually be used to optimize the selection and dosing of medications for individual patients. The discovery and translation of inherited determinants of drug response and somatically acquired genetic variants in cancer are prominent pharmacogenomic components of these and other initiatives.

Diagnostic testing

Before being used in a clinical setting, genetic tests must meet certain criteria concerning their analytical validity, clinical validity and clinical utility²⁰. Developing genetic tests with analytical validity is nontrivial for pharmacogenes such as *CYP2D6* (ref. 21) because the gene is subject to germline copy number variation and the formation of hybrid gene fusions that are difficult to assay and interpret reliably. Clinical utility involves assessing whether the use of the test leads to improved health outcomes for patients who are subject to testing, as well as an assessment of the risks that occur as a result of testing. There is substantial difference of opinion as to precisely which outcomes constitute clinical utility^{22,23}. Some studies have broadened their scope to assess the impact of testing on the entire health-care system, including comparing the costs of genetic testing with those of other health-care interventions, as well as understanding how such testing can influence the behaviour of clinicians. For example, Hong Kong introduced a policy to screen for the *HLA-B*1502* allele before prescribing antiepileptic drugs because patients with the allele are at high risk of developing severe skin reactions to the drug carbamazepine. But the policy led clinicians to forego prescribing carbamazepine at all and instead they began to prescribe phenytoin. Phenytoin can also cause severe skin reactions, but the risk factors are not well defined, as for carbamazepine, so the overall incidence of severe skin reactions remained unchanged²⁴. For the purposes of this Review, we focus on the clinical utility of pharmacogenomic testing for individual patients without considering the possible public-health consequences of changes in prescribing behaviour.

As the cost of sequencing continues to fall, many have predicted that every individual will in the not-too-distant future have their germline genome sequenced early in life and the results will be available for clinical use throughout their lifetime. Assuming that this prediction is realized (to some extent, at least), we call for a shift away from the debate

over whether patients should be tested for specific pharmacogenes before they are prescribed specific drugs, and instead, we suggest moving towards a model in which clinicians are provided with guidelines on how to interpret and deploy genetic variants to improve their prescribing. This assumption underlies the efforts of the Clinical Pharmacogenetics Implementation Consortium (CPIC)^{25,26}, an open, international non-profit group that creates standardized guidelines on how to use genomic data to inform prescribing. The guidelines are evidence-based, peer-reviewed and publicly available.

A number of factors are used to determine whether there is enough evidence to support the analytical validity, clinical validity and clinical utility of a pharmacogenomic test and to warrant its use in guiding the prescription of medications²². The analytical validity will depend on the quality of the data from genetic tests as well as the test’s performance characteristics, such as the positive and negative predictive values. Many types of data can be used to evaluate clinical validity and utility, including the penetrance of genetic variation on drug effects, which can be determined from retrospective studies. The mechanism or mechanisms by which genetic variation influences drug effects or a relevant endophenotype (an intermediate phenotype, such as drug-metabolizing enzyme activity) can also be used. Additionally, data can be gathered from *in vivo* pharmacokinetic or other functional studies, *in vitro* functional studies,

Table 1 | Actionable germline genetic variation and associated medications

Genetic variation	Medications
<i>TPMT</i>	Mercaptopurine, thioguanine, azathioprine
<i>CYP2D6</i>	Codeine, tramadol, tricyclic antidepressants
<i>CYP2C19</i>	Tricyclic antidepressants, clopidogrel, voriconazole
<i>VKORC1</i>	Warfarin
<i>CYP2C9</i>	Warfarin, phenytoin
<i>HLA-B</i>	Allopurinol, carbamazepine, abacavir, phenytoin
<i>CFTR</i>	Ivacaftor
<i>DPYD</i>	Fluorouracil, capecitabine, tegafur
<i>G6PD</i>	Rasburicase
<i>UGT1A1</i>	Irinotecan, atazanavir
<i>SLCO1B1</i>	Simvastatin
<i>IFNL3 (IL28B)</i>	Interferon
<i>CYP3A5</i>	Tacrolimus

From ref. 44 (accessed on 7 May 2015). See ref. 44 for updates.

and preclinical and clinical studies that link pharmacological effects or drug concentrations to genetic variation. Further sources of data include case reports, family studies, and randomized clinical trials that compare the outcomes of genetics-based prescribing versus prescribing that is not based on genetic-test results. Other factors that are considered when deciding on the actionability of pharmacogenomic variation include the therapeutic index of a drug (the ratio of the toxic dose to the therapeutic dose), the severity of the drug's toxicity, the severity of the underlying disease and the consequences of suboptimal prescribing.

An important consideration for the actionability of a gene–drug relationship is the availability of an alternative therapy, which may partly depend on the mechanism of the gene–drug association. If the gene affects the pharmacokinetics of a drug (such as in the *CYP3A5*-mediated catabolism of the immunosuppressant drug tacrolimus), there could be substantial published evidence to support a dose adjustment in much the same way as doses are often adjusted according to age or kidney or liver function. Such dose-adjustment decisions are particularly defensible if the drug is one for which therapeutic drug monitoring (based on the concentration in the blood) is readily available. When tests indicate that a drug is either unlikely to be effective or could have unacceptable adverse effects in patients with a particular genotype, the recommendation for an alternative therapy will depend on the balance of evidence for both the efficacy and possible toxicity of the alternative. For example, individuals who are homozygous for inactive *CYP2D6* alleles are unable to convert the analgesic drug codeine into its active metabolite, morphine, but can respond to several other opiate analgesics²⁷. If genetic tests indicate an extremely high risk for a serious adverse event — for example, carriers of the *HLA-B*57:01* allele have a high risk of hypersensitivity to the antiretroviral drug abacavir^{28,29} — the alternative therapy should be equally effective but have an acceptable risk of adverse effects (that may or may not be influenced by other genetic variants).

Some treatment decisions are less clear cut. For example, there are substantial data to show that patients with two defective *CYP2D6* alleles are more likely to experience recurrence of breast cancer after treatment with tamoxifen. This is because these patients produce much lower levels of the active tamoxifen metabolite, endoxifen, than those without defective alleles^{30–35}. However, it is unclear whether the best alternative is another drug (a different selective oestrogen-receptor modulator, for instance) or an altered dose of tamoxifen, particularly in premenopausal women for whom there is a shortage of data to support alternative treatments. Such cases are the most difficult to resolve: although clear from pharmacogenomic testing that the drug or drug dose is suboptimal in a patient with the high-risk genotype, a lack of clinical data for alternative therapies makes it difficult to recommend other medications.

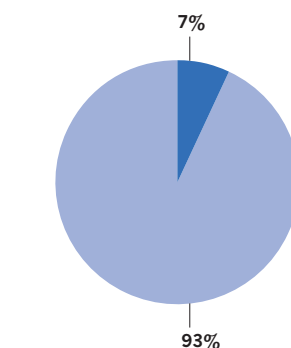
The CPIC considers all such evidence when deciding which gene–drug pairs are clinically actionable. Given the high bar for clinical actionability, the number of actionable genes — those with at least one actionable high-risk diplotype — is small, and the list of medications for which clinical actions are recommended (pharmacogenetically high-risk drugs) is relatively short (Table 1). There are also additional medications that include pharmacogenomic information in their labels^{36–38}. However, not all are actionable. Sometimes, information on genetic variation is included when the effects are modest and therefore do not require changes to be made to the prescribing section of the drug label. This information has also been included for some drug labels despite weak or conflicting evidence.

At present, there are only two examples of actionable pharmacogenes that also carry a disease risk: the gene *UGT1A1* and Gilbert's syndrome³⁹, and the gene *G6PD* and haemolytic anaemia⁴⁰. Thus, many of the ethics concerns that affect the clinical implementation of disease-risk genomics have less relevance for pharmacogenomics⁴¹.

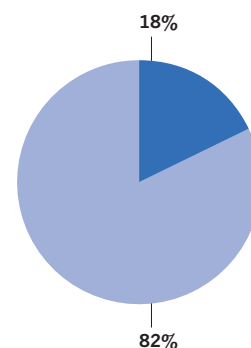
Clinical implementation of pharmacogenomics

Genetic variants that influence the clinical effects of some medications can now be reliably assayed in the clinical setting. Prescribing decisions for such clinically actionable gene–drug pairs should be influenced by these genetic-test results.

FDA-approved medications
(*n* = 1,200)



Prescriptions in the United States
(*n* = 4 billion)



■ Affected by actionable pharmacogenes
■ Not affected by actionable pharmacogenes

Figure 2 | Medications affected by actionable pharmacogenes.

Approximately 7% of FDA-approved medications are affected by actionable inherited pharmacogenes (left), and approximately 18% of US outpatient prescriptions are affected by actionable germline pharmacogenes (right)⁴⁵, which demonstrates that several pharmacogenetically high-risk drugs are commonly prescribed.

Drugs and genes

More than 1,200 individual molecular entities have been approved as drugs by the US Food and Drug Administration (FDA)⁴², the European Medicines Agency (EMA)³⁶ or by Japan's Pharmaceuticals and Medical Devices Agency (PMDA)³⁸. Although about 15% of the medications approved by the FDA and EMA contain pharmacogenomic information on their label^{36,43}, only a subset of the corresponding pharmacogenes is deemed actionable. As summarized in Fig. 2, 7% of medications have actionable germline pharmacogenetics. These correspond to CPIC level A or B gene–drug pairs⁴⁴ for which genetic information should or could be used to change the prescription pattern of the relevant drug. In the United States, these medications constitute 18% of all prescriptions, which indicates that pharmacogenomically high-risk medications are slightly overrepresented in highly prescribed medications (Fig. 2)⁴⁵. So far, only 16 of the roughly 19,000 human genes are considered to be clinically actionable for germline pharmacogenomics⁴⁴. Most human germline genetic variation is unlikely to be actionable for the prescription of medication, and pharmacogenomics is unlikely to be useful for improving the prescription of the majority of drugs. However, for the relatively small set of medications for which genomics is actionable, prescribing could be optimized if genetic testing were more widely and appropriately deployed in the clinic. In the meantime, the number of such actionable gene–drug pairs continues to grow, albeit at a slow pace.

Somatically acquired genomic variation

Genetic variants that are specific to cancer tissue represent a special case of pharmacogenomics. Somatic variation can identify which types of malignancy are likely to respond to various anticancer agents^{46,47}. The recognition that cancer tissue can be distinguished from normal tissue by the presence of specific genomic abnormalities pre-dates the Human Genome Project. For example, in the 1980s and 1990s ploidy in neuroblastoma⁴⁸ and cytogenetic abnormalities in acute lymphoblastic leukaemia⁴⁹ were used to determine the composition and strength of the cytotoxic chemotherapy required to treat these cancers. The genetic testing of malignancies has since become more specific in response to the development of anticancer agents that are directed at — or more effective at — treating tumours that harbour certain acquired genetic variants (Table 2). Although the FDA generally requires that diagnostics be developed alongside targeted anticancer agents, the EMA is less stringent⁵⁰. However, proposed changes to the European Union's framework would lead to greater harmonization between the United States and Europe⁵¹.

Table 2 | Actionable somatic genetic variants in cancer cells and associated medications

Genetic abnormality	HGVS nomenclature*	Target†	Medications	Disease
<i>AKT</i> mut (act)	p.Glu17Lys	mTOR	Sirolimus, everolimus	RCC
<i>BCR-ABL</i> (SV)	t(9;22)(q34.1;q11.21)	ABL	Imatinib, dasatinib	CML, Ph ⁺ ALL
<i>BCR-ABL</i> (SV, mut)	p.Val299Leu	ABL	Bosutinib, nilotinib	Imatinib-resistant CML
<i>BCR-ABL</i> (T135I)	p.Thr135Ile	ABL	Ponatinib	CML, Ph ⁺ ALL
<i>BCR-ABL</i> (SV)	t(9;22)(q34.1;q11.21)	SRC	Dasatinib	CML, Ph ⁺ ALL
<i>BRCA1/2</i> variants	Too numerous to list	PARP	Olaparib	Ovarian cancer
<i>BRAF</i> SNVs (V600E/K)	p.Val600Glu, p.Val600Lys, p.Val600Asp	BRAF	Dabrafenib, vemurafenib	Melanoma
<i>BRAF</i> SNVs (V600)	p.Val600Glu, p.Val600Lys, p.Val600Asp	MEK	Trametinib	Melanoma
<i>EGFR</i> (ex 19 del, SNV L858R)	p.Glu746_Ala750del, p.Leu858Arg	EGFR	Afatinib, erlotinib	NSCLC (EGFR ⁺)
<i>EGFR</i> mut (act, amp)	p.Glu746_Ala750del, p.Leu858Arg	EGFR	Gefitinib	NSCLC (EGFR ⁺)
<i>EGFR</i> ⁺ and WT <i>KRAS</i>	N/A	EGFR	Cetuximab, panitumumab	EGFR ⁺ colon cancer (WT <i>KRAS</i>)
<i>EML-ALK</i> (SV)	inv(2)(p21p23)	ALK	Crizotinib	NSCLC
<i>FLT3</i> CNV (amp)	p.D600_L601insFREYED, p.Asp835Tyr	FTL3	Sunitinib, sorafenib	AML
<i>HER2</i> (amp)	N/A	ERBB2	Lapatinib, trastuzumab	HER2 ⁺ breast cancer
<i>KIT</i> mut (act)	p.Trp557_Lys558del, p.Asp579del, p.Val559Asp	KIT	Imatinib, sunitinib	RCC, GIST
<i>PDGFR</i> (mut or SV)	p.Asp842Val	PDGFR	Sunitinib, imatinib	RCC, GIST, pancreatic cancer
<i>PI3K</i> (mut or amp)	PIK3CA p.Glu542Lys, p.Glu545Lys; p.His1047Arg, p.His1047Leu	PI3K	Idelalisib	CLL, NHL
<i>RARA</i> (SV, gene fusion)	t(15;17)(q24;q21)	RARA	Tretinoin, alitretinoin	APL, CTCL, Kaposi sarcoma
<i>RARA</i> (SV, gene fusion)	t(15;17)(q24;q21)	RARA	Arsenic trioxide	APL
<i>SMO</i> mut (act)	p.Trp535Leu, p.Arg199Trp, p.Arg562Gln	Smoothen	Vismodegib	Basal cell carcinoma
<i>VHL</i> (mut)	Too numerous to list	VEGFR	Sorafenib	RCC, hepatic cancer, thyroid cancer
<i>VEGF</i> (mut)	N/A	VEGF	Ziv-aflibercept	Colon cancer

Medications targeting normal cell surface proteins that are expressed on some tumour cells (for example, ER, PR, CD20, CD30, CD52) are not included in this summary of drugs that target proteins with aberrant expression or function due to somatic genetic variants. Act, activating; ALCL, anaplastic large cell lymphoma; ALL, acute lymphoblastic leukaemia; AML, acute myeloid leukaemia; amp, amplification (typically by CNV); APL, acute promyelocytic leukaemia; CLL, chronic lymphocytic leukaemia; CML, chronic myeloid leukaemia; CNV, copy number variant; CTCL, cutaneous T-cell lymphoma; del, deletion; ER, oestrogen receptor; ex, exon; GIST, gastrointestinal stromal tumour; HGVS, Human Genome Variation Society; ins, insertion; mut, mutation; N/A, not applicable; NHL, non-Hodgkin lymphoma; NSCLC, non-small-cell lung cancer; Ph, Philadelphia chromosome; PR, progesterone receptor; RCC, renal cell carcinoma; SNV, single nucleotide variant; SV, structural variant; WT, wild type.

*Only representative examples of known mutations are shown.

†In general, targets are protein products encoded by the gene listed.

Reactive and pre-emptive testing

Initially, clinical testing was deployed through single-gene pharmacogenetic tests — a practice that evolved from the coupling of strong ‘monogenic’ gene–drug associations with limitations in genotyping technology⁵². In this model, genetic tests are ordered one at a time on a reactive basis: if the patient might need a pharmacogenetically high-risk drug, the clinician orders the applicable test. However, improvements in technology mean that it is now possible to interrogate multiple genes in a single assay at lower cost than for multiple single-gene tests.

Most human diseases, including cancer, are influenced by multiple genes and genetic variants. Likewise, the pharmacokinetics and pharmacological effects of most medications are determined by multiple gene products, such as drug-metabolizing enzymes, drug-transporting membrane proteins, drug targets and disease-modifying genes. Many of the actionable pharmacogenes that have been identified to date exert a strong effect on the pharmacokinetics or pharmacodynamics of their associated drug, which make them easy to identify. The strong effects of the polymorphism in the gene *TPMT* on the risk of haematopoietic toxicity from thiopurine medications are an excellent illustration of how ‘low-hanging fruit’ are often merely the first step down a polygenic path.

Subsequent studies showed that when the dosage of mercaptopurine is adjusted in response to *TPMT* testing, polymorphisms in other genes — such as *ITPA* — begin to surface as important⁵³. Furthermore, polymorphisms in other genes that form part of the same pharmacological pathway can emerge as being significant for populations of a different ancestry. An example of this is the strong influence of an inherited variant found in the gene *NUDT15* on thiopurine toxicity. *NUDT15* variants are extremely uncommon in individuals of European and African ancestry, but are relatively common in people of Asian descent⁵⁴. This explains the high frequency of thiopurine intolerance in Asian populations despite the low frequency of *TPMT* variants. When a genome-wide association study of thiopurine intolerance was performed in a population that comprised people of European, Asian, African and Native American ancestry, germline variants in both *TPMT* and *NUDT15* reached genome-wide significance for their association with thiopurine intolerance⁵⁵. *TPMT* variants were revealed to be the major determinant of the tolerated dose of thiopurine medications in patients of European and African ancestry, whereas *NUDT15* was the major genetic determinant in patients of Asian and Native American ancestry. Moreover, the metabolism and effects of anticancer agents, including thiopurines,

can be affected by both germline and somatic genetic variation⁵⁶, which further increases the complexity of cancer pharmacogenomics.

Several other examples exist in which more than one gene is clinically actionable for a given medication. These include the anticoagulant warfarin, which is affected by both *CYP2C9* and *VKORC1* (ref. 19), and tricyclic antidepressants, which are affected by both *CYP2C19* and *CYP2D6*. Given that a single gene can affect more than one medication (Table 1), there are potential benefits of genotyping a panel of pharmacogenomic variants that apply to a number of drugs that a patient might receive in their lifetime. There is increasing evidence to show that genotyping multiple genes in a single assay is more cost-effective, uses the DNA in the sample more efficiently, and facilitates the pre-emptive availability of genetic-test information. Such multigene panels can change practice from a reactive approach (in which a fresh genetic test is ordered every time it is required) to a pre-emptive approach (in which a single sample is assessed for many likely-to-be actionable genes at the same time), thereby providing the patient with a lifetime's worth of test results. Several groups have already begun to implement pre-emptive multigene panels for pharmacogenomics^{57–61}, but the practice is by no means widespread at present.

Clinical implementation

A number of barriers prevent the widespread use of pre-emptive multigene panels to guide the prescription of drugs. These include the lack of incentives for clinicians to conduct tests or implement procedures that might prevent adverse events. There are relatively few studies that prove the cost-effectiveness of pharmacogenetic testing⁶². Although a multigene panel approach is less expensive than ordering tests for one pharmacogene at a time, there are no data to assess the cost-effectiveness of the panel approach when implemented early on in life and used throughout a patient's lifetime. Many health-care systems do not provide financial reimbursement for preventive-medicine services or for pre-emptive screening, which creates a barrier to pharmacogenetic testing in the clinic^{63,64}.

The cost and complexity of the computational approaches needed to identify, catalogue, prioritize and interpret genetic variants that influence prescribing decisions present another barrier to the clinical uptake of pharmacogenomics testing. Although an increasing number of computational tools for analysing genetic variation are becoming available, a substantial level of expertise and manual interpretation is still required to apply this analysis successfully in the clinic (Fig. 3). Computational tools for CDS, triggered by patient-specific alerts, prompt and guide clinicians to use genetic information when prescribing affected drugs^{45,65,66}. The costs associated with the use of pharmacogenomics in clinical practice are quickly shifting away from laboratory testing towards the process of linking genetic-test results with evidence-based decisions that will robustly guide the prescription of medications and will be updated as the latest evidence emerges. However, it is unclear who should take responsibility for updating and paying for such interpretations.

A further barrier to the clinical uptake of pharmacogenomic testing is the lack of clear guidelines for translating genetic variation into actionable recommendations. Professional societies and other guideline-generating groups sometimes disagree on whether to proceed with pharmacogenetic testing and, if so, how. Examples for which there has been disagreement include the drugs warfarin⁶⁷ and clopidogrel⁶⁸. A common reason for the lack of support for genetic testing is the dearth of randomized prospective controlled trials that compare genetically guided testing with conventional therapy. Also, many professional societies and guideline-generating groups approach their evaluation of pharmacogenomic tests from the standpoint of whether the clinician is obligated to order the genetic test^{52,67–69}. However, as inexpensive multigene tests become available, the question is shifting from whether to order a genetic test to how existing genetic-test results can and should be used to influence prescribing decisions. For inherited genetic variation, the CPIC has undertaken the task of creating guidelines that focus on how genetic-test results should be translated into specific prescribing

actions. The Royal Dutch Pharmacists Association took a similar approach^{70,71}. Multiple resources exist to guide the selection of cancer drugs on the basis of somatically acquired genetic variation (Table 2), although these constantly change as new evidence arises^{72–75}.

As deep sequencing becomes more widespread, further variants will be discovered in pharmacogenes⁷⁶. The challenge will be to catalogue and annotate these variants. Given the importance of rare variants for both inherited²⁴ and cancer-related pharmacogenes, publicly available and easily updatable resources such as PharmGKB, ClinGen and ClinVar will be essential for providing the computational CDS in health-care record systems with up-to-date recommendations that are based on genetic-test results^{77–79}. Current heterogeneity in genetic-variation databases and health-care record systems, coupled with a lack of a common ontology, limits interoperability and hinders the use of pharmacogenetic test results longitudinally as well as across each of the health-care systems the patient must navigate. Several groups are working to standardize the language of pharmacogenetic testing^{80–85}, with the aim of creating terminology that can drive CDS across health-care record systems. Initiatives such as an Institute of Medicine Roundtable on Translating Genomic-Based Research for Health and the CPIC are working to create terms and language that can be directly uploaded into

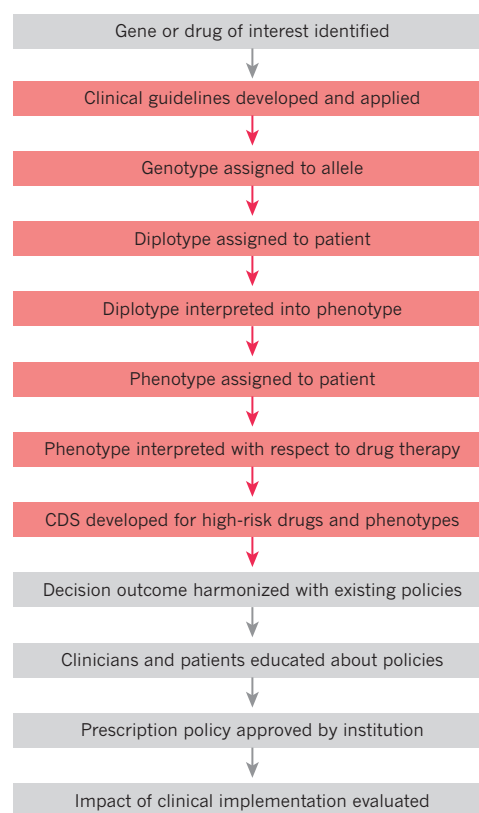


Figure 3 | Bringing pharmacogenomic testing to the clinic. Multiple steps are required before pharmacogenomic testing can be used in the clinic. First, genes or drugs for actionability are identified. Guidelines for all actionable inherited pharmacogenes, such as those developed by CPIC²⁵, exist for subsequent steps that are shaded in red. Next, genotypes are assigned to alleles, and diplotypes are assigned to patients. Diplotypes are then translated to phenotypes (that describe gene function), followed by their interpretation with respect to drug therapy. Appropriate CDS is deployed to provide clinicians with recommendations for prescribing medications. Pharmacogenetic considerations are harmonized with other policies for affected medications, including therapeutic drug monitoring, if applicable. Clinicians and patients are educated as to what the results mean, what type of CDS and information to expect and which medication might be affected. Institutional oversight committees can then approve prescribing recommendations and policies, if desired. Last, clinical and prescribing outcomes are audited by various groups to evaluate the impact of the clinical implementation of pharmacogenomics.

the CDS of electronic health records. However, heterogeneity across systems will initially slow the creation and uptake of CDS that facilitates the use of pharmacogenetic information⁸⁴.

Other barriers include the insatiable desire for more evidence, a lack of education in clinicians, the scarcity of evidence-based implementation systems, and concerns about incidental or secondary findings from genetic testing, not to mention inertia in health-care systems^{69,86–89}. These barriers are not unique to pharmacogenomics, and the energy required to overcome them is likely to come from multiple sources that range from the ‘push’ of patient advocates to the ‘pull’ of the courtroom. The resistance by some clinicians to the ordering or use of pharmacogenetic testing can be perplexing to patients and their caregivers. For example, an advocate for paediatric patients expressed a disquieting lay perspective: “I am mystified by the resistance to a simple blood test that might save children’s lives.”¹⁷ As the general public becomes more aware of the potential for genetic tests to improve the prescription of medications, including through direct-to-consumer testing^{90,91}, its advocacy could grow even stronger. In Hawaii, the attorney general brought a lawsuit against the manufacturer of clopidogrel because it marketed its drug in the state without warning that a high percentage of the Hawaiian population has inherited low-function alleles of the gene *CYP2C19*, which encodes the enzyme required to convert clopidogrel to its active metabolite⁹². The case asserts that it was already known that *CYP2C19* variant allele frequencies are higher in East Asian and Pacific Island populations, which comprise about 40% and 10% of the Hawaiian population, respectively. It also states that there was abundant evidence to show that the antithrombotic effects of clopidogrel are diminished in patients with low *CYP2C19* activity (predisposing them to an increased incidence of cardiovascular events such as stent thrombosis)^{93,94}. From an educational perspective, multiple accrediting agencies are calling for pharmacogenomics to become part of the curricula for health-care students, trainees and advanced practitioners⁹⁵. The availability of pharmacogenomic educational tools continues to grow^{96–98}. Although early adopters of clinical pharmacogenomics are now establishing methods to advance the treatment of patients, broad clinical implementation remains elusive.

Facilitating clinical use

Around the world, many groups are working to share resources that will facilitate the clinical implementation of germline pharmacogenetic tests⁹⁹. The European Pharmacogenetics Implementation Consortium is an international body whose goal is to improve therapy by integrating pharmacogenetic information into clinical care¹⁰⁰. The Royal Dutch Pharmacists Association has also made efforts to facilitate implementation^{70,71}. In the United States, members of the NIH Pharmacogenomics Research Network organized the Translational Pharmacogenetics Project^{57–61,101–103}, which is dedicated to sharing best practices for the clinical implementation of CPIC pharmacogenomics guidelines. Meanwhile, the Electronic Medical Records and Genomics (eMERGE) and Implementing Genomics in Practice (IGNITE) networks are testing pharmacogenetic implementation strategies^{88,104,106}. In Thailand and Singapore, where the *HLA-B*15:02* variant is widespread and strongly predisposes patients to severe skin toxicity after treatment with specific drugs, pharmacogenetic testing is common^{106–108}. The Genomic Medicine Alliance¹⁰⁹ facilitates the clinical use of pharmacogenomics and has created a database that links drugs with genes¹¹⁰. Population admixture between previously isolated, diverse human populations must also be considered as part of the global effort towards clinical implementation¹¹¹.

Future directions

Clinicians are accustomed to making prescribing decisions on the basis of patient characteristics such as age, kidney or liver function, drug–drug interactions and personal preferences. Much of this takes place without optimal CDS to assist in compiling these characteristics and matching them with evidenced-based choices on medications

and doses. As CDS improves and becomes more widespread, and as the evidence to support pharmacogenomic testing continues to grow, momentum for the clinical implementation of pharmacogenomics should accelerate. Going forward, there is a growing body of evidence to suggest that pharmacogenomics will become an important component of evidence-based precision medicine. ■

Received 7 May; accepted 7 August 2015.

- Hughes, D. & Andersson, D. I. Evolutionary consequences of drug resistance: shared principles across diverse targets and organisms. *Nature Rev. Genet.* **16**, 459–471 (2015).
- Evans, W. E. & Relling, M. V. Moving towards individualized medicine with pharmacogenomics. *Nature* **429**, 464–468 (2004).
- A review of pharmacogenomics, from discovery to the clinic.**
- Alving, A. S., Carson, P. E., Flanagan, C. L. & Ickes, C. E. Enzymatic deficiency in primaquine-sensitive erythrocytes. *Science* **124**, 484–485 (1956).
- Roll-Hansen, N. The crucial experiment of Wilhelm Johannsen. *Biol. Phil.* **4**, 303–329 (1989).
- Motulsky, A. G. Drug reactions enzymes, and biochemical genetics. *J. Am. Med. Assoc.* **165**, 835–837 (1957).
- Kalow, W. & Genest, K. A method for the detection of atypical forms of human serum cholinesterase; determination of dibucaine numbers. *Can. J. Biochem. Physiol.* **35**, 339–346 (1957).
- Vogel, F. Moderne problem der humangenetik. *Ergeb. Inn. Med. U. Kinderheik.* **12**, 52–125 (1959).
- Vesell, E. S. & Page, J. G. Genetic control of drug levels in man: antipyrine. *Science* **161**, 72–73 (1968).
- Hughes, H. B., Biehl, J. P., Jones, A. P. & Schmidt, L. H. Metabolism of isoniazid in man as related to the occurrence of peripheral neuritis. *Am. Rev. Tuberc.* **70**, 266–273 (1954).
- Price Evans, D. A., Manley, K. A. & McKusick, V. A. Genetic control of isoniazid metabolism in man. *Br. Med. J.* **2**, 485–491 (1960).
- Blum, M., Demierre, A., Grant, D. M., Heim, M. & Meyer, U. A. Molecular mechanism of slow acetylation of drugs and carcinogens in humans. *Proc. Natl Acad. Sci. USA* **88**, 5237–5241 (1991).
- Vatsis, K. P., Martell, K. J. & Weber, W. W. Diverse point mutations in the human gene for polymorphic *N*-acetyltransferase. *Proc. Natl Acad. Sci. USA* **88**, 6333–6337 (1991).
- Gonzalez, F. J. et al. Characterization of the common genetic defect in humans deficient in debrisoquine metabolism. *Nature* **331**, 442–446 (1988).
- Ingelman-Sundberg, M. Pharmacogenomic biomarkers for prediction of severe adverse drug reactions. *N. Engl. J. Med.* **358**, 637–639 (2008).
- Wang, L., McLeod, H. L. & Weinshilboum, R. M. Genomics and drug response. *N. Engl. J. Med.* **364**, 1144–1153 (2011).
- Yates, C. R. et al. Molecular diagnosis of thiopurine S-methyltransferase deficiency: genetic basis for azathioprine and mercaptopurine intolerance. *Ann. Intern. Med.* **126**, 608–614 (1997).
- Marshall, E. Preventing toxicity with a gene test. *Science* **302**, 588–590 (2003).
- Carr, D. F., Alfirevic, A. & Pirmohamed, M. Pharmacogenomics: Current state-of-the-art. *Genes (Basel)* **5**, 430–443 (2014).
- Pirmohamed, M., Kamali, F., Daly, A. K. & Wadelius, M. Oral anticoagulation: a critique of recent advances and controversies. *Trends Pharmacol. Sci.* **36**, 153–163 (2015).
- A discussion of how to evaluate the benefits of individualized therapy, and how population differences can complicate this, for warfarin — one of the most important clinically actionable drugs.**
- Burke, W. in *Current Protocols in Human Genetics* Unit 9.15, 9.15.1–9.15.7 (John Wiley & Sons, 2009).
- Gaedigk, A. Complexities of *CYP2D6* gene analysis and interpretation. *Int. Rev. Psychiatry* **25**, 534–553 (2013).
- Grosse, S. D. & Khoury, M. J. What is the clinical utility of genetic testing? *Genet. Med.* **8**, 448–450 (2006).
- Scott, S. A. Personalizing medicine with clinical pharmacogenetics. *Genet. Med.* **13**, 987–995 (2011).
- Chen, Z., Liew, D. & Kwan, P. Effects of a *HLA-B*15:02* screening policy on antiepileptic drug use and severe skin reactions. *Neurology* **83**, 2077–2084 (2014).
- Relling, M. V. & Klein, T. E. CPIC: Clinical Pharmacogenetics Implementation Consortium of the Pharmacogenomics Research Network. *Clin. Pharmacol. Ther.* **89**, 464–467 (2011).
- Caudle, K. E. et al. Incorporation of pharmacogenomics into routine clinical practice: the Clinical Pharmacogenetics Implementation Consortium (CPIC) guideline development process. *Curr. Drug Metab.* **15**, 209–217 (2014).
- Crews, K. R. et al. Clinical Pharmacogenetics Implementation Consortium guidelines for cytochrome P450 2D6 genotype and codeine therapy: 2014 update. *Clin. Pharmacol. Ther.* **95**, 376–382 (2014).
- Mallal, S. et al. *HLA-B*57:01* screening for hypersensitivity to abacavir. *N. Engl. J. Med.* **358**, 568–579 (2008).
- Martin, M. A. et al. Clinical Pharmacogenetics Implementation Consortium guidelines for *HLA-B* genotype and abacavir dosing: 2014 update. *Clin. Pharmacol. Ther.* **95**, 499–500 (2014).
- Province, M. A., Altman, R. B. & Klein, T. E. Interpreting the *CYP2D6* results from the International Tamoxifen Pharmacogenetics Consortium. *Clin. Pharmacol. Ther.* **96**, 144–146 (2014).

31. Ratain, M. J., Nakamura, Y. & Cox, N. J. *CYP2D6* genotype and tamoxifen activity: understanding interstudy variability in methodological quality. *Clin. Pharmacol. Ther.* **94**, 185–187 (2013).
32. Brauch, H. *et al.* Tamoxifen use in postmenopausal breast cancer: *CYP2D6* matters. *J. Clin. Oncol.* **31**, 176–180 (2013).
33. Irvin, W. J. Jr *et al.* Genotype-guided tamoxifen dosing increases active metabolite exposure in women with reduced *CYP2D6* metabolism: a multicenter study. *J. Clin. Oncol.* **29**, 3232–3239 (2011).
34. Brauch, H. & Schwab, M. Prediction of tamoxifen outcome by genetic variation of *CYP2D6* in post-menopausal women with early breast cancer. *Br. J. Clin. Pharmacol.* **77**, 695–703 (2014).
35. Province, M. A. *et al.* *CYP2D6* genotype and adjuvant tamoxifen: meta-analysis of heterogeneous study populations. *Clin. Pharmacol. Ther.* **95**, 216–227 (2014).
36. Ehmann, F. *et al.* Pharmacogenomic information in drug labels: European Medicines Agency perspective. *Pharmacogenomics J.* **15**, 201–210 (2015).
37. Tutton, R. Pharmacogenomic biomarkers in drug labels: what do they tell us? *Pharmacogenomics J.* **15**, 297–304 (2014).
38. Ishiguro, A., Yagi, S. & Uyama, Y. Characteristics of pharmacogenomics/biomarker-guided clinical trials for regulatory approval of anti-cancer drugs in Japan. *J. Hum. Genet.* **58**, 313–316 (2013).
39. Bosma, P. J. *et al.* The genetic basis of the reduced expression of bilirubin UDP-glucuronosyltransferase 1 in Gilbert's syndrome. *N. Engl. J. Med.* **333**, 1171–1175 (1995).
40. Beutler, E. G6PD deficiency. *Blood* **84**, 3613–3636 (1994).
41. McCarthy J. J., McLeod, H. L. & Ginsburg, G. S. Genomic medicine: a decade of successes, challenges, and opportunities. *Sci. Transl. Med.* **5**, 189sr4 (2013).
42. Kinch, M. S., Haynesworth, A., Kinch, S. L. & Hoyer, D. An overview of FDA-approved new molecular entities: 1827–2013. *Drug Discov. Today* **19**, 1033–1039 (2014).
43. US Food and Drug Administration. Table of Pharmacogenomic Biomarkers in Drug Labeling. *US Food and Drug Administration* <http://www.fda.gov/drugs/scienceresearch/researchareas/pharmacogenetics/ucm083378.htm> (2015).
44. PharmGKB. CPIC Genes/Drugs. *PharmGKB* <https://www.pharmgkb.org/cpic/pairs> (2015).
45. Dunnenberger, H. M. *et al.* Preemptive clinical pharmacogenetics implementation: current programs in five US medical centers. *Annu. Rev. Pharmacol. Toxicol.* **55**, 89–106 (2015).
- A detailed discussion of a pre-emptive approach to the clinical implementation of pharmacogenetics that contains quantitative information on the use of medications that are subject to genetic actionability.**
46. Wheeler, H. E., Maitland, M. L., Dolan, M. E., Cox, N. J. & Ratain, M. J. Cancer pharmacogenomics: strategies and challenges. *Nature Rev. Genet.* **14**, 23–34 (2013).
47. McLeod, H. L. Cancer pharmacogenomics: early promise, but concerted effort needed. *Science* **339**, 1563–1566 (2013).
- An analysis of the considerations for use of both somatic and inherited germline genetic lesions in prescription of anticancer drugs.**
48. Look, A. T., Hayes, F. A., Nitschke, R., McWilliams, N. B. & Green, A. A. Cellular DNA content as a predictor of response to chemotherapy in infants with unresectable neuroblastoma. *N. Engl. J. Med.* **311**, 231–235 (1984).
49. Pui, C. H., Crist, W. M. & Look, A. T. Biology and clinical significance of cytogenetic abnormalities in childhood acute lymphoblastic leukemia. *Blood* **76**, 1449–1463 (1990).
50. Senderowicz, A. M. & Pfaff, O. Similarities and differences in the oncology drug approval process between FDA and European Union with emphasis on *in vitro* companion diagnostics. *Clin. Cancer Res.* **20**, 1445–1452 (2014).
51. Pignatti, F. *et al.* Cancer drug development and the evolving regulatory framework for companion diagnostics in the European Union. *Clin. Cancer Res.* **20**, 1458–1468 (2014).
52. Flockhart, D. A., Skaar, T., Berlin, D. S., Klein, T. E. & Nguyen, A. T. Clinically available pharmacogenomics tests. *Clin. Pharmacol. Ther.* **86**, 109–113 (2009).
53. Stocco, G. *et al.* Genetic polymorphism of inosine triphosphatase pyrophosphatase is a determinant of mercaptopurine metabolism and toxicity during treatment for acute lymphoblastic leukemia. *Clin. Pharmacol. Ther.* **85**, 164–172 (2009).
54. Yang, S. K. *et al.* A common missense variant in *NUDT15* confers susceptibility to thiopurine-induced leukopenia. *Nature Genet.* **46**, 1017–1020 (2014).
55. Yang, J. J. *et al.* Inherited *NUDT15* variant is a genetic determinant of mercaptopurine intolerance in children with acute lymphoblastic leukemia. *J. Clin. Oncol.* **33**, 1235–1242 (2015).
56. Cheng, Q. *et al.* Karyotypic abnormalities create discordance of germline genotype and cancer cell phenotypes. *Nature Genet.* **37**, 878–882 (2005).
57. Bielinski, S. J. *et al.* Preemptive genotyping for personalized medicine: design of the right drug, right dose, right time—using genomic data to individualize treatment protocol. *Mayo Clin. Proc.* **89**, 25–33 (2014).
58. Gottesman, O. *et al.* The CLIPMERGE PGx Program: clinical implementation of personalized medicine through electronic health records and genomics-pharmacogenomics. *Clin. Pharmacol. Ther.* **94**, 214–217 (2013).
59. Fernandez, C. A. *et al.* Concordance of DMET plus genotyping results with those of orthogonal genotyping methods. *Clin. Pharmacol. Ther.* **92**, 360–365 (2012).
60. Johnson, J. A. *et al.* Implementing personalized medicine: development of a cost-effective customized pharmacogenetics genotyping array. *Clin. Pharmacol. Ther.* **92**, 437–439 (2012).
61. Oetjens, M. T. *et al.* Assessment of a pharmacogenomic marker panel in a polypharmacy population identified from electronic medical records. *Pharmacogenomics J.* **14**, 735–744 (2013).
62. Buchanan, J., Wordsworth, S. & Schuh, A. Issues surrounding the health economic evaluation of genomic technologies. *Pharmacogenomics J.* **14**, 1833–1847 (2013).
63. Schroeder, S. A. & Frist, W. Phasing out fee-for-service payment. *N. Engl. J. Med.* **368**, 2029–2032 (2013).
64. Levy, K. D. *et al.* Prerequisites to implementing a pharmacogenomics program in a large health-care system. *Clin. Pharmacol. Ther.* **96**, 307–309 (2014).
65. Overby, C. L. *et al.* Physician attitudes toward adopting genome-guided prescribing through clinical decision support. *J. Pers. Med.* **4**, 35–49 (2014).
66. Crawford, D. C. *et al.* eMERGEing progress in genomics—the first seven years. *Front. Genet.* **5**, 184 (2014).
67. Cavallari, L. H. & Nutescu, E. A. Warfarin pharmacogenetics: to genotype or not to genotype, that is the question. *Clin. Pharmacol. Ther.* **96**, 22–24 (2014).
68. Chan, N. C. *et al.* Role of phenotypic and genetic testing in managing clopidogrel therapy. *Blood* **124**, 689–699 (2014).
69. Stanek, E. J. *et al.* Adoption of pharmacogenomic testing by US physicians: results of a nationwide survey. *Clin. Pharmacol. Ther.* **91**, 450–458 (2012).
70. Swen, J. J. *et al.* Pharmacogenetics: from bench to byte. *Clin. Pharmacol. Ther.* **83**, 781–787 (2008).
71. Swen, J. J. *et al.* Pharmacogenetics: from bench to byte—an update of guidelines. *Clin. Pharmacol. Ther.* **89**, 662–673 (2011).
- A survey of clinically actionable germline genetic variants and affected medications, with basic prescribing advice.**
72. Yeh, P. *et al.* DNA-mutation inventory to refine and enhance cancer treatment (DIRECT): a catalog of clinically relevant cancer mutations to enable genome-directed anticancer therapy. *Clin. Cancer Res.* **19**, 1894–1901 (2013).
73. Van Allen, E. M., Wagle, N. & Levy, M. A. Clinical analysis and interpretation of cancer genome data. *J. Clin. Oncol.* **31**, 1825–1833 (2013).
- An overview of databases that can be used to match somatic cancer-specific genetic variants with targeted anticancer drugs.**
74. Abrams, J. *et al.* in *2014 American Society of Clinical Oncology Education Book 71–76* (American Society of Clinical Oncology, 2014).
75. Agúndez, J. A., Esguevilas, G., Amo, G. & García-Martín, E. Clinical practice guidelines for translating pharmacogenomic knowledge to bedside. Focus on anticancer drugs. *Front. Pharmacol.* **5**, 188 (2014).
76. Gordon, A. S. *et al.* Quantifying rare, deleterious variation in 12 human cytochrome P450 drug-metabolism genes in a large-scale exome dataset. *Hum. Mol. Genet.* **23**, 1957–1963 (2014).
77. Landrum, M. J. *et al.* ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* **42**, D980–D985 (2014).
78. Rehm, H. L. *et al.* ClinGen—the Clinical Genome Resource. *N. Engl. J. Med.* **372**, 2235–2242 (2015).
79. Whirl-Carrillo, M. *et al.* Pharmacogenomics knowledge for personalized medicine. *Clin. Pharmacol. Ther.* **92**, 414–417 (2012).
80. Percha, B., & Altman, R. B. Inferring the semantic relationships of words within an ontology using random indexing: applications to pharmacogenomics. *AMIA Annu. Symp. Proc.* **2013**, 1123–1132 (2013).
81. Samwald, M. & Freimuth, R. R. Making data on essential pharmacogenes available for every patient everywhere: the Medicine Safety Code initiative. *Pharmacogenomics J.* **14**, 1529–1531 (2013).
82. Bell, G. C. *et al.* Development and use of active clinical decision support for preemptive pharmacogenomics. *J. Am. Med. Assoc.* **21**, e93–e99 (2014).
83. Zhu, Q. *et al.* Harmonization and semantic annotation of data dictionaries from the Pharmacogenomics Research Network: a case study. *J. Biomed. Inform.* **46**, 286–293 (2013).
84. Overby, C. L. *et al.* Opportunities for genomic clinical decision support interventions. *Genet. Med.* **15**, 817–823 (2013).
85. Miñarro-Giménez, J. A., Blagec, K., Boyce, R. D., Adlansnig, K. P. & Samwald, M. An ontology-based, mobile-optimized system for pharmacogenomic decision support at the point-of-care. *PLoS ONE* **9**, e93769 (2014).
86. Haga, S. B. *et al.* Survey of genetic counselors and clinical geneticists' use and attitudes toward pharmacogenetic testing. *Clin. Genet.* **82**, 115–120 (2012).
87. Haga, S. B., Burke, W., Ginsburg, G. S., Mills, R. & Agans, R. Primary care physicians' knowledge of and experience with pharmacogenetic testing. *Clin. Genet.* **82**, 388–394 (2012).
88. Weber, G. M., Mandl, K. D. & Kohane, I. S. Finding the missing link for big biomedical data. *J. Am. Med. Assoc.* **311**, 2479–2480 (2014).
89. Hayden, E. C. Geneticists push for global data-sharing. *Nature* **498**, 16–17 (2013).
90. Prainsack, B. & Vayena, E. Beyond the clinic: 'direct-to-consumer' genomic profiling services and pharmacogenomics. *Pharmacogenomics J.* **14**, 403–412 (2013).
91. Caulfield, T. DTC genetic testing: pendulum swings and policy paradoxes. *Clin. Genet.* **81**, 4–6 (2012).
92. State of Hawaii Department of the Attorney General. *State of Hawaii News Release 2014-09* <http://ag.hawaii.gov/wp-content/uploads/2014/01/News-Release-2014-09.pdf> (2014).
93. Shuldiner, A. R. *et al.* Association of cytochrome P450 2C19 genotype with the antiplatelet effect and clinical efficacy of clopidogrel therapy. *J. Am. Med. Assoc.* **302**, 849–857 (2009).
- Primary findings that showed that *CYP2C19* genetic variation affects the effectiveness and bleeding risk from clopidogrel, and led to an FDA 'black box warning' for the drug label.**
94. Mega, J. L. *et al.* Reduced-function *CYP2C19* genotype and risk of adverse clinical outcomes among patients treated with clopidogrel predominantly for PCI: a meta-analysis. *J. Am. Med. Assoc.* **304**, 1821–1830 (2010).

95. Manolio, T. A. & Green, E. D. Leading the way to genomic medicine. *Am. J. Med. Genet. C. Semin. Med. Genet.* **166C**, 1–7 (2014).
96. Manolio, T. A., Murray, M. F. & Inter-Society Coordinating Committee for Practitioner Education in Genomics. The growing role of professional societies in educating clinicians in genomics. *Genet. Med.* **16**, 571–572 (2014).
97. Korf, B. R. *et al.* Framework for development of physician competencies in genomic medicine: report of the Competencies Working Group of the Inter-Society Coordinating Committee for Physician Education in Genomics. *Genet. Med.* **16**, 804–809 (2014).
98. Wiener, C. M., Thomas, P. A., Goodspeed, E., Valle, D. & Nichols, D. G. “Genes to society”—the logic and process of the new curriculum for the Johns Hopkins University School of Medicine. *Acad. Med.* **85**, 498–506 (2010).
99. Manolio, T. A. *et al.* Global implementation of genomic medicine: we are not alone. *Sci. Transl. Med.* **7**, 290ps13 (2015).
100. Becquemont, L. *et al.* Practical recommendations for pharmacogenomics-based prescription: 2010 ESF-UB Conference on Pharmacogenetics and Pharmacogenomics. *Pharmacogenomics* **12**, 113–124 (2011).
101. Shuldiner, A. R. *et al.* The Pharmacogenomics Research Network Translational Pharmacogenetics Program: overcoming challenges of real-world implementation. *Clin. Pharmacol. Ther.* **94**, 207–210 (2013).
102. Pulley, J. M. *et al.* Operational implementation of prospective genotyping for personalized medicine: the design of the Vanderbilt PREDICT project. *Clin. Pharmacol. Ther.* **92**, 87–95 (2012).
- A description of the benefits and efficiency of implementing a large and innovative pre-emptive pharmacogenetics programme at a major medical centre.**
103. O'Donnell, P. H. *et al.* Adoption of a clinical pharmacogenomics implementation program during outpatient care—initial results of the University of Chicago “1,200 Patients Project”. *Am. J. Med. Genet. C. Semin. Med. Genet.* **166C**, 68–75 (2014).
104. Gottesman, O. *et al.* The Electronic Medical Records and Genomics (eMERGE) Network: past, present, and future. *Genet. Med.* **15**, 761–771 (2013).
105. Rasmussen-Torvik, L. J. *et al.* Design and anticipated outcomes of the eMERGE-PGx project: a multicenter pilot for preemptive pharmacogenomics in electronic health record systems. *Clin. Pharmacol. Ther.* **96**, 482–489 (2014).
106. Rattanaipapong, W., Koopitakajorn, T., Praditsithikorn, N., Mahasirimongkol, S. & Teerawattananon, Y. Economic evaluation of HLA-B*15:02 screening for carbamazepine-induced severe adverse drug reactions in Thailand. *Epilepsia* **54**, 1628–1638 (2013).
107. Toh, D. S. *et al.* Building pharmacogenetics into a pharmacovigilance program in Singapore: using serious skin rash as a pilot study. *Pharmacogenomics J.* **14**, 316–321 (2014).
108. Sukasem, C., Puangpetch, A., Medhasi, S. & Tassaneeyakul, W. Pharmacogenomics of drug-induced hypersensitivity reactions: challenges, opportunities and clinical implementation. *Asian Pac. J. Allergy Immunol.* **32**, 111–123 (2014).
109. Cooper, D. N. *et al.* Bridging genomics research between developed and developing countries: the Genomic Medicine Alliance. *Pers. Med.* **11**, 615–623 (2014).
110. Dalabira, E. *et al.* DruGeVar: an online resource triangulating drugs with genes and genomic biomarkers for clinical pharmacogenomics. *Public Health Genomics* **17**, 265–271 (2014).
111. Bonifaz-Peña, V. *et al.* Exploring the distribution of genetic markers of pharmacogenomics relevance in Brazilian and Mexican populations. *PLoS ONE* **9**, e112640 (2014).
112. Relling, M. V. *et al.* Clinical Pharmacogenetics Implementation Consortium (CPIC) guidelines for rasburicase therapy in the context of G6PD deficiency genotype. *Clin. Pharmacol. Ther.* **96**, 169–174 (2014).

Acknowledgements The authors thank their many colleagues who have collaborated in research to elucidate the pharmacogenomics of medications used to treat their patients, as well as those who have collaborated to translate pharmacogenomics to optimize the care of patients. They also thank the many patients and parents who have willingly participated in their research. This work was supported in part by NIH grants P50 GM115279, R01 CA36401, R24 GM115264, R01 CA142665, P30 CA21765 and by the American Lebanese Syrian Associated Charities.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare competing financial interests: see go.nature.com/4muiva for details. Readers are welcome to comment on the online version of this paper at go.nature.com/4muiva. Correspondence should be addressed to M.V.R. (mary.relling@stjude.org).

Gene therapy returns to centre stage

Luigi Naldini^{1,2}

Recent clinical trials of gene therapy have shown remarkable therapeutic benefits and an excellent safety record. They provide evidence for the long-sought promise of gene therapy to deliver ‘cures’ for some otherwise terminal or severely disabling conditions. Behind these advances lie improved vector designs that enable the safe delivery of therapeutic genes to specific cells. Technologies for editing genes and correcting inherited mutations, the engagement of stem cells to regenerate tissues and the effective exploitation of powerful immune responses to fight cancer are also contributing to the revitalization of gene therapy.

Gene therapy has long fascinated scientists, clinicians and the general public because of its potential to treat a disease at its genetic roots. This is achieved by counteracting or replacing a malfunctioning gene within the cells adversely affected by the condition. As simple as the concept sounds, the hurdles to put it into practice are daunting. Gene transfer must overcome complex cellular and tissue barriers to deliver new genetic information into the target cell to drive proficient expression of a therapeutic molecule without disrupting essential regulatory mechanisms. The gene-corrected cells must be present in large enough quantities to reverse the condition, escape immunological recognition and survive in the long term, or be able to transmit the modification to their progeny to sustain the benefit. Several gene-therapy trials have been performed in the past two decades for inherited diseases, cancer and chronic infections, but only a few reported clear clinical benefits and in some, individuals experienced severe adverse events related to the vectors. Overall, concern and scepticism rose over the further deployment of these strategies. But these attitudes are radically changing. A number of phase I/II gene-therapy clinical trials have reported remarkable evidence of efficacy and safety for the treatment of various severe inherited diseases of the blood, immune and nervous systems, including primary immunodeficiencies, leukodystrophies, thalassaemia, haemophilia and retinal dystrophy, as well as cancers such as B-cell malignancies (Table 1). All of these trials exploit improved vector technologies to deliver therapeutic genes. In some trials, genetic material is transferred into haematopoietic stem cells (HSCs) or T lymphocytes (T cells) *ex vivo*, and in others hepatocytes in the liver or photoreceptors in the retina are targeted directly *in vivo*. Here I review the most relevant clinical results, highlight progress in gene-transfer technologies and in our understanding of the biological processes that underpin these advances and discuss the challenges and outlook for gene therapy.

Haematopoietic-stem-cell gene therapy

HSCs have long been a preferred target for *ex vivo* gene therapy¹. Genetic modification of self-maintaining multipotent HSCs would ensure a steady supply of their gene-corrected progeny in the body. These cells have the potential to treat conditions that manifest when mature haematopoietic lineages fail to develop or to function correctly. Given the self-renewing nature of HSCs and the need to ensure that genetic modifications are passed on to their progeny, gene correction must be stably

introduced into cellular chromatin, either by vector-mediated transgene insertion or by *in situ* gene editing.

Past trials of HSC gene therapy report clear benefits to people with selected conditions, which proved the therapeutic potential of the strategy^{2–10}. However, they also highlight the limitations and risks of using early generation vectors that are based on gammaretroviruses (γ -RVs)^{11–15}. For example, these techniques offered only a limited ability to transfer genes into the most primitive progenitor cells, giving rise to the low-level and transient appearance of gene-corrected haematopoietic cells *in vivo*. Leukaemia related to vector insertion near oncogenes also occurred in a fraction of patients during their long-term follow-up. The development of vectors with improved efficacy and safety in preclinical models, such as lentiviral vectors, has renewed interest in the approach¹.

HSC gene therapies that incorporate these new vectors have been tested in severe inherited diseases of the immune system (Wiskott–Aldrich syndrome (WAS) and X-linked severe combined immunodeficiency (SCID-X1))^{16–18} and blood (β -thalassaemia)¹⁹, and in neurodegenerative storage diseases (adrenoleukodystrophy and metachromatic leukodystrophy)^{20–22} (see Table 1). In children with SCID-X1, the development and function of the immune system is impaired owing to deficiencies in the receptors for certain cytokines that are essential for immune-cell development, whereas those with WAS have deficiencies in a cytoskeletal adaptor that is required for assembling the immunological synapse²³. Consequently, these children succumb to infection or, in the case of WAS, haemorrhage because of an accompanying platelet deficiency. People with β -thalassaemia major — the most severe form of β -thalassaemia — fail to express the haemoglobin- β chain, which leads to ineffective erythropoiesis and severe anaemia that requires frequent blood transfusions and an iron-chelation regimen. Children with early onset adrenoleukodystrophy or metachromatic leukodystrophy are affected by defective myelination and consequently experience glial- and neural-cell degeneration in the central and peripheral nervous systems. They are unable to break down some of the metabolites of myelin because of a deficiency in the peroxisomal ATP-binding cassette transporter in adrenoleukodystrophy²⁴, or in the lysosomal enzyme arylsulphatase A in metachromatic leukodystrophy²⁵. Such patients undergo rapid and irreversible deterioration of their motor, sensory and cognitive performance, which leads to death within a few years. HSC transplantation, in which HSCs are transferred from a healthy donor to a recipient with a specific condition, can virtually cure SCID-X1, WAS and β -thalassaemia²⁶. HSC transplantation can also

¹San Raffaele Telethon Institute for Gene Therapy (TIGET), San Raffaele Scientific Institute, 20132 Milan, Italy. ²Vita Salute San Raffaele University, 20132 Milan, Italy.

Table 1 | Gene-therapy clinical trials highlighted in this Review

Disease	Vector and strategy	Number of patients*	Follow-up (months)	Patient status and biological and clinical outcomes*	Clinical-trial identifier	References
HSC-based gene therapy						
Wiskott–Aldrich syndrome	Lentiviral vector; <i>ex vivo</i> gene transfer into CD34 ⁺ cells.	7	10 to 60	All patients AAW; stable engraftment with transduced cells; persistent clinical benefit and safety.	NCT01515462	16 and L.N.†
Wiskott–Aldrich syndrome	Lentiviral vector; <i>ex vivo</i> gene transfer into CD34 ⁺ cells.	7	9 to 42	6 patients AAW, 1 patient died of a pre-existing infection; stable engraftment with transduced cells; persistent clinical benefit and safety.	NCT01347242 NCT01347346 NCT02333760	17
X-linked severe combined immunodeficiency	Self-inactivating γ -RV; 9 <i>ex vivo</i> gene transfer into CD34 ⁺ cells.	9	12 to 39	8 patients AAW, 1 patient died of an infection; stable engraftment with transduced cells; persistent clinical benefit and safety in 7 patients; 1 patient failed to engraft and underwent HSC transplantation.	NCT01410019 NCT01175239 NCT01129544	18
β -Thalassaemia major	Lentiviral vector; <i>ex vivo</i> gene transfer into CD34 ⁺ cells.	3‡	24 to 72	2 patients stably engrafted with transduced cells, 1 patient became transfusion independent; 1 patient failed to engraft and received rescue cells.	N/A	19, M. Cavazzana and BlueBird Bio†
β -Thalassaemia major	Lentiviral vector; <i>ex vivo</i> gene transfer into CD34 ⁺ cells.	2‡	15	Stable engraftment with transduced cells; transfusion independence and safety in both patients.	NCT02151526	M. Cavazzana and BlueBird Bio†
β -Thalassaemia major	Lentiviral vector; <i>ex vivo</i> gene transfer into CD34 ⁺ cells.	5§	1 to 6	Stable engraftment with transduced cells; safety and transfusion independence in the first 2 evaluable patients.	NCT01745120	BlueBird Bio†
Adrenoleukodystrophy	Lentiviral vector; <i>ex vivo</i> gene transfer into CD34 ⁺ cells.	4	54 to 101	Stable engraftment with transduced cells and safety in all patients; persistent clinical benefit in 3 patients.	N/A	20, 21 and P. Aubourg†
Metachromatic leukodystrophy	Lentiviral vector; <i>ex vivo</i> gene transfer into CD34 ⁺ cells.	20	3 to 60	Stable engraftment with transduced cells and safety in all patients; persistent clinical benefit in all late-infantile patients who were treated when presymptomatic.	NCT01560182	22 and L.N.†
Liver-directed gene therapy						
Haemophilia B	AAV8 vector; intravenous administration.	10	16 to 48	No inhibitors; persistent FIX expression; in high-dose group, mean FIX levels of $5.1 \pm 1.7\%$ seen in all 6 treated patients.	NCT00979238	40
Haemophilia B	AAV8 vector; intravenous administration.	7	Up to 12	No inhibitors; persistent FIX expression in 1 patient.	NCT01687608	110
T-cell immunotherapy for cancer						
B-cell lymphoma or CLL	γ -RV; <i>ex vivo</i> gene transfer into T cells; CAR-modified anti-CD19 cells.	15	1 to 23	8 CRs, 4 PRs; ORR 80%.	NCT00924326	51
B-cell ALL	γ -RV; <i>ex vivo</i> gene transfer into T cells; CAR-modified anti-CD19 cells.	5	1 to 4	5 CRs; ORR 100%; 4 patients subsequently underwent allo-HSC transplantation as per clinical-study design, 1 patient was ineligible for HSC transplantation and relapsed.	NCT01044069	52
B-cell ALL	Lentiviral vector; <i>ex vivo</i> gene transfer into T cells; CAR-modified anti-CD19 cells.	30	1 to 24¶	27 CRs; ORR 90%; 19 patients remained in remission, 3 of these patients underwent HSC transplantation; 7 patients relapsed, 3 of these relapses occurred after loss of transduced T cells.	NCT01626495 NCT01029366	53
B-cell ALL or lymphoma	γ -RV; <i>ex vivo</i> gene transfer into T cells; CAR-modified anti-CD19 cells.	21	Median = 10#	14 CRs (at day 28); ORR 67%; 10 patients subsequently underwent HSC transplantation.	NCT01593696	54
Retinal gene therapy						
Type 2 Leber congenital amaurosis	AAV2 vector; unilateral subretinal administration.	5	36	In all 5 patients, stable improvement in visual sensitivity seen.	NCT00516477	94, 97
Type 2 Leber congenital amaurosis	AAV2 vector; unilateral subretinal administration.	3	54 to 72	In all 3 patients, improvement in visual sensitivity seen at 6 months, which increased for 1 to 3 years and then declined.	NCT00481546	95
Type 2 Leber congenital amaurosis	AAV2 vector; unilateral subretinal administration.	12	36	In 6 patients, improvement in visual sensitivity seen, which peaked at 6 to 12 months and then declined.	NCT00643747	96

AAW, alive and well; ALL, acute lymphocytic leukaemia; CLL, chronic lymphocytic leukaemia; CR, complete response; N/A, not applicable; ORR, overall response rate; PR, partial response.

*As stated in the referenced publications or updated by personal communication. †Personal communication. ‡ β^0/β^E genotype. §2 β^0/β^E , 2 β^0/β^0 , 1 β^0/β^+ genotypes. ¶Median = 7. #51 days to HSC transplantation median.

arrest the progression of leukodystrophies if performed before or near the time of onset, although the outcome of the procedure is more satisfactory in adrenoleukodystrophy²⁴ than in metachromatic leukodystrophy²⁷. Allogeneic HSC transplantation confers a considerable risk of morbidity and death, particularly when performed between human leukocyte antigen (HLA)-mismatched individuals. Even in the few cases in which an HLA-matched family donor can be found, the risk of morbidity remains substantial²⁸. HSC gene therapy can address this unmet medical need, especially when no matched HSC donor is available.

HSC gene therapy is administered by *ex vivo* gene transfer into haematopoietic progenitors. First, the cells are purified using the CD34 surface marker from other leukocytes harvested from the bone marrow or mobilized peripheral blood of the recipient. Next, they are cultured for 2–4 days in the presence of growth-stimulating cytokines while being exposed to vectors carrying an expression cassette for the corrective transgene. Before the modified cells can be administered, the recipient is treated with a pre-conditioning chemotherapy regimen. This depletes endogenous progenitors and differentiated cells in the bone marrow — as well as the lymphoid organs, in some cases — and favours engraftment of the *ex vivo* gene-corrected cells. Preconditioning results in considerable early morbidity owing to transitory blood-cell depletion, immunodeficiency and mucosal damage, which place the recipient at risk of severe infection. It also causes delayed morbidity owing to the risk of developing chemotherapy-induced secondary tumours and infertility²⁸. Several HSC gene-therapy trials have attempted to alleviate the morbidity associated with preconditioning by lowering the dosage and combination of chemotherapeutic drugs that are administered in comparison to the regimens used in conventional HSC transplantation. However, the impact of changing these drug regimens on the risks and benefits of the therapy is yet to be determined in broader comparative studies and through long-term patient follow-up. Some preconditioning regimens also deplete endogenous microglia progenitors, which live in the central nervous system (CNS). In patients with leukodystrophies who have been treated with HSC gene therapy, these cells are replaced with their gene-corrected counterparts, which migrate to CNS tissues and produce functional progeny that can clear stored myelin metabolites. In the case of metachromatic leukodystrophy, these modified cells also release functional lysosomal enzyme for the cross-correction of other tissue-resident cells^{29–31}.

All HSC gene-therapy trials performed with lentiviral vectors report stable and high-level reconstitution of haematopoiesis with gene-corrected cells in most recipients, with the most recent trials observing up to 90% reconstitution in some patients^{16,17,19–22}. Because all of the haematopoietic lineages tested (myelo-monocytic, megakaryocytic, erythroid, natural killer (NK), and B lymphoid and T lymphoid) contain gene-corrected cells and appear to receive a steady input of newly formed modified cells — most evident in the short-lived myeloid lineages — gene transfer must have occurred in self-renewing, multipotent progenitors that effectively engrafted to the recipients. The extent of haematopoietic reconstitution with gene-corrected cells varies among the patients and trials. This reflects the impact of the disease and its correction on the different haematopoietic lineages as well as the potency of each vector batch and the type of preconditioning regimen administered. Remarkably, most patients gain substantial benefits from gene therapy that can even exceed those observed after successful HSC transplantation from an allogeneic donor. This is especially true for a metachromatic leukodystrophy trial²² based on the rationale that gene transfer could drive higher lysosomal enzyme expression in haematopoietic cells than in normal cells, which would enable these cells to provide an increased supply of enzyme while homing in to affected tissues. The strategy allowed the unsatisfactory outcome of HSC transplantation to be overcome in metachromatic leukodystrophy. Despite an increasing number of patients being treated in these HSC gene-therapy trials (see Table 1), there has been no report of adverse events related to lentiviral vectors. At the most recent follow-up of patients included in these trials, the first of whom was treated almost seven years ago in the adrenoleukodystrophy trial²⁰, those who had successfully engrafted with modified cells were reported to be in a good

condition — that is, they are either disease free or their disease has been stable or in remission since the time of treatment — and are leading near-normal lives. These individuals would probably have already succumbed to their disease if untreated with gene therapy^{24,25}.

Liver-directed gene therapy

The liver has long been a preferred target for *in vivo* gene therapy³². This major internal organ and central metabolic hub receives an abundant blood supply through an extensive bed of sinusoids with highly permeable walls — a structure that facilitates the access of blood-borne particles, such as viruses, to hepatocytes. Hepatocytes are long-lived and robust protein factories that can efficiently release their products into the blood circulation. Stable transgene delivery to the liver could therefore provide a strategy for treating several inherited metabolic diseases and plasma-protein deficiencies, notably those of coagulation factors. Major hurdles to liver-directed gene therapy include the potential toxicity of an acute inflammatory response to the bolus administration of viral particles into the bloodstream, and the inactivation of these particles by pre-existing virus-specific antibodies and clearance by phagocytes that line the sinusoidal walls of the liver and spleen. If humoral or cellular immunity is triggered against the transgene product, the therapeutic activity of gene therapy could be inhibited in the circulation or the modified cells might be eliminated by cytotoxic T cells, respectively^{33,34}.

Liver-directed gene therapy has been tested mainly in the treatment of severe haemophilia B using vectors derived from the human parvovirus, adeno-associated virus (AAV)³⁵. On intravascular injection, some AAV-vector serotypes effectively target hepatocytes and remain stably associated with the modified cells, mostly as nonintegrated episomes within the nucleus. In animal models, this delivery strategy does not induce an immune response against the transgene product and instead favours the development of transgene-specific tolerance³². The persistence of AAV-vector genomes that contain an expression cassette for coagulation factor IX (FIX) — a protein that is absent or defective in haemophilia B — within nonproliferating hepatocytes can give rise to a stable supply of functional FIX into the bloodstream³⁶. FIX expression as low as 1% of the normal level can turn severe haemophilia B into a milder form of disease. This partial reconstitution alleviates the risk of spontaneous haemorrhages, which are detrimental to the joints and potentially lethal if they occur in the brain. It can also reduce the need for intravenous prophylactic factor-replacement therapy, a treatment that imposes a substantial burden on patients and is only available at high cost and in countries with well-developed health-care systems³⁷. Although earlier trials demonstrated the safety of delivering AAV vectors into the human bloodstream, they reported only limited and transient FIX expression. This is probably because of low levels of gene transduction in the liver and a delayed cellular immune response that targeted AAV components persisting within modified hepatocytes and triggered their elimination^{38,39}. More recently, a trial that used a new AAV serotype (AAV8) and an improved cassette design was able to overcome these limitations to demonstrate vector-dose-dependent FIX expression of up to 6% of the normal level after a single, well-tolerated vector infusion in a peripheral vein⁴⁰. Three years later, FIX expression remained stable in most patients, although some had to be given a transient immunosuppressive corticosteroid treatment at the first sign of hepatocellular injury in the 4–8 weeks following treatment. Most treated patients were able to reduce or abrogate the need for prophylactic factor replacement, which substantially improved their quality of life.

T-cell immunotherapy for cancer

T cells are also popular targets for *ex vivo* gene therapy. Such therapies aim mostly at boosting the adaptive immune response against cancer and chronic infections such as HIV^{41,42}. Autologous T cells can be harvested readily from the peripheral blood and expanded *ex vivo*. Cells are then transduced with a γ -RV or lentiviral vector expressing an exogenous T-cell antigen receptor (TCR) that is specific to a cancer-associated antigen or an antiviral molecule, and infused back into the patient. The use of T cells for cancer immunotherapy arose from seminal observations of objective

BOX 1

Lentiviral vectors

Major hurdles for haematopoietic-stem-cell (HSC) gene therapy include achieving efficient *ex vivo* gene transfer into long-term repopulating HSCs, preventing activation of oncogenes by the nearby integration of a vector and controlling transgene expression to avoid ectopic or constitutive expression that leads to toxicity¹. As compared to early generation gammaretroviral vectors (γ-RVs), HIV-derived lentiviral vectors result in more efficient gene transfer and stable, robust transgene expression in HSCs and their multilineage progeny. Extensive preclinical work indicated important features in vector biology and design that affect genotoxicity and highlighted strategies to alleviate it^{111–117}. The self-inactivating long terminal repeats (LTRs) and integration-site preferences of lentiviral vectors were shown to substantially alleviate insertional genotoxicity. When tested in γ-RVs, the self-inactivating LTR design was shown to improve the safety of this platform as well¹⁸. Retrospective analysis of several earlier trials suggests that disease background, transgene function, *ex vivo* culture and the efficiency of host repopulation can all influence the likelihood that insertional genotoxicity will manifest in a trial^{13,14}. These data helped to shape the ideas that not all integrating vectors have the same effects and that genome-wide integration of improved vector designs, although still a mutagenic event, can be tolerated in the absence of aggravating circumstances¹¹⁸.

Self-inactivating lentiviral vectors are also being used to engineer T cells with chimaeric antigen receptors (CARs) or T-cell antigen receptors for use in adoptive immunotherapy for the treatment of cancer. The advantages of this new platform in comparison to earlier-generation γ-RVs, which perform satisfactorily in this cell target, are yet to be fully established. Lentiviral vectors are thought to give rise to more robust and stable transgene expression in T cells *in vivo*, and could facilitate more efficient and versatile *ex vivo* gene transfer while supporting coordinated expression of multiple transgenes^{41,42,55,119}. These advantages will become more relevant as the gene-therapy field implements refined strategies, such as improved T-cell manipulation to preserve T memory stem cells^{59–61}, or more demanding cell-engineering tasks, such as the co-expression of multiple CARs (to improve specificity) or a conditional safety switch/suicide gene (to improve safety)¹²⁰.

tumour responses, which sometimes led to complete tumour regression, after the infusion of *ex vivo*-expanded autologous tumour-infiltrating lymphocytes in patients with advanced melanoma⁴¹. The underlying mechanisms of these responses have been traced to the *ex vivo* activation and amplification — and subsequent *in vivo* persistence — of tumour-specific cytotoxic T cells that already exist in small numbers within some tumours but are suppressed by the local microenvironment. These data support the therapeutic potential of an adaptive immune response against some tumours when it is released from endogenous suppressive signals and the reactive T cells are infused at high numbers, which leads to a favourable effector–target ratio. Preconditioning patients with a lymphoid-depleting regimen before infusion of the *ex vivo*-expanded T cells improves their engraftment and activity *in vivo*⁴³.

The genetic transfer of an exogenous TCR to a T cell can bypass the need to find preexisting cancer specificities within the patient's T-cell population. It also allows the exploitation of synthetic high-affinity TCRs, which might otherwise be purged *in vivo* by thymic selection, to redirect the specificity of autologous cells. Ideally, such TCRs would target tumour-associated antigens from endogenous proteins that are not being expressed — at the same time or at detectable levels — by normal tissues, such as carcinoembryonic or cancer testis antigens. Alternatively, they could target tumour-driver mutations that are commonly found

within a certain type of tumour and that cannot be lost, even to evade immune clearance. Recent studies, however, indicate that most spontaneous or elicited tumour-specific immune reactivity is instead directed against 'passenger' neoantigens, which uniquely originate from random mutations that accumulate within individual tumours^{44–47}. This type of immunotherapy must therefore be highly personalized and will require both the identification of candidate neoantigens from the tumour exome or proteome and the retrieval or generation of the cognate TCR recognition sequences for *ad hoc* T-cell genetic engineering.

Early clinical trials of T-cell gene transfer with TCRs directed against tumour-associated antigens reported partial tumour responses and, in some instances, off-tumour reactivity that led to tissue damage and severe adverse events^{48–50}. More recently, gene transfer was used to introduce synthetic chimaeric antigen receptors (CARs) to T cells. CARs combine the binding specificity of an antibody against a cancer-associated surface marker with one or more intracellular signalling domains from the TCR and costimulatory receptor complexes⁴². The strategy holds a number of advantages over the use of conventional TCRs. For instance, antigen recognition is not restricted by HLA, and antibody specificities previously validated for safe and specific cancer recognition *in vivo* can be exploited by the incorporation of single-chain antibody derivatives into the CARs. Furthermore, these engineered T cells can be fully activated when meeting their target. Trials that deployed CARs directed against the B-cell surface molecule CD19 reported dramatic benefits in patients with B-cell malignancies, who experienced durable clinical responses, including complete remission, with mostly manageable toxicity^{51–54}. These results have spurred enormous interest in further developing this approach⁵⁵.

Reasons for success in recent clinical trials

The positive clinical outcomes discussed in this Review provide long-sought evidence for the elusive promise of gene therapy to deliver lasting therapeutic benefit, or even a 'cure', for an otherwise terminal or severely disabling condition after a single treatment. These results could represent the rewarding outcome of the effort to rationally improve vectors in the laboratory, as discussed in Box 1 and Box 2. Advances in vector manufacturing and characterization, leading to batches with higher potency and greater purity, could also be facilitating the increased transduction of target cells with lower adverse effects, even when using the same vector design as in earlier trials. A better understanding of the biological processes and specific cell types that are involved in the success or failure of each gene therapy has also helped to improve methods for *ex vivo* cell handling and led to more effective monitoring and management of patients. Indeed, increasing confidence in the positive outcome of recent trials is now built on both clinical observation and advanced molecular readouts that provide evidence for the efficacy and safety of gene therapy.

In recent lentiviral-vector-based HSC gene-therapy trials, vector insertional analysis was used to track individual clones in the reconstituted haematopoiesis and showed that polyclonal reconstitution by transduced stem cells had taken place without the emergence of dominant clones whose behaviour could be attributed to gain-of-function vector insertional mutagenesis. An in-depth molecular analysis of the reconstituted haematopoiesis now being performed in different diseases, including several that were treated with different types of vector, is facilitating the first reliable comparative assessment of vector-induced events in patients^{16,17,20–22,56}. The emerging picture illustrates that the absence of adverse clinical events in recent trials is accompanied by a remarkably different landscape of vector-insertion distribution in patients. In the earlier trials, which were performed with long terminal repeat (LTR)-competent γ-RV vectors, there was frequent and early generation of dominant clones, whose expansion appeared to be driven by the altered expression of cancer genes that were targeted by vector insertion^{11–13,14,57}. However, this finding is rarely observed when newer vectors, such as lentiviral vectors, with self-inactivating LTRs and different insertion-site preferences are used^{16,17,19–22}. In parallel, retrospective reconstruction of multistep leukaemogenesis from the earlier trials supports the view that

vector insertion provided a first mutagenic hit that potentially led to the development of leukaemia. Progression of the disease was facilitated by accumulation of further vector-independent mutations in the expanding clone, which were favoured by concurrent precipitating factors such as oligoclonal reconstitution, selective pressure for a survival or growth advantage, transgene toxicity or stressed haematopoiesis resulting from the underlying disease^{12,13,14}. Overall, these findings support the idea that improved efficiency of HSC transduction and transplantation, together with improved vector design, can substantially alleviate, but not eliminate, concerns about genotoxicity. Once it has been established that vector insertion is mostly neutral to cell behaviour, tracking clonal activity in the reconstituted haematopoiesis can highlight important biological features of haematopoietic regeneration. These include the pattern of activity and stability of the transduced HSCs, which provides the foundation for predicting long-term maintenance of the therapeutic benefit^{16,22}. Tracking clones can also help to establish the reliability of current models of lineage relationship and haematopoietic hierarchy based on xenotransplantation studies⁵⁸.

In adoptive T-cell therapy, a subpopulation of T memory stem cells plays an important part in supporting long-term efficacy. This knowledge arose from an improved understanding and phenotypic characterization of T-cell differentiation pathways, and the identification of the cells that are responsible for the sustained generation of effector activity *in vivo*⁵⁹. Optimization of *ex vivo* cell culture⁶⁰ and *in vivo* tracking of gene-marked T cells⁶¹ also contributed to this discovery, which should facilitate the design of more effective trials. Epitope mapping and tracking of cancer-targeting T cells in patients who are undergoing adoptive T-cell therapies and checkpoint blockade drug treatment provide direct evidence for the potential to evoke cancer-specific effector T cells in some types of cancer, which could then be exploited to clear advanced metastatic disease after *ex vivo* amplification or genetic editing^{44–47}.

In liver-directed gene therapy, detailed epitope-specific immune monitoring of antiviral T-cell responses on *in vivo* AAV gene delivery is able to account partially for earlier failures to observe long-term stable FIX expression^{34,39}. This also helped to uncover strategies for controlling a delayed cytotoxic adaptive response to viral components that persist within transduced cells, which include administering transient immunosuppressive drugs at the first occurrence of such a response⁴⁰ and engineering viral capsids to bypass dominant responses^{62–64}.

The enhanced efficacy observed in recent gene-therapy trials could also reflect improved study designs. These have been made possible by a deeper understanding of vector–host interactions and by the application of more rational patient-selection criteria, which are more easily adopted after the initial ‘proof of safety’ of administration by vector in humans has been obtained. Improved trial designs can enhance efficacy in several ways. For example, treating early symptomatic (or even presymptomatic) individuals instead of people with an advanced stage of a disease provides a greater chance for the therapy to work before tissue damage has become irreversible. Vector doses can be set at levels that approach therapeutic efficacy instead of minimal biological activity, as is usually done to alleviate risks when establishing the safety of a treatment in phase I/II trials. However, this cautionary approach often denies trial participants the chance to gain immediate benefit from the therapy, as well as any future benefit (because they become immune to the vector). By administering more aggressive preconditioning regimens, space can be made for the gene-corrected cells. Vector-induced adverse events can also be neutralized by preemptive or prompt pharmacological treatment, such as administering anti-inflammatory drugs or steroids to suppress innate or adaptive immune responses, respectively.

Gene-therapy trials often address rare diseases for which little is known about the natural history and genotype–phenotype relationships. Enriching such knowledge helps in the design of more effective trials because it becomes possible to recruit the patients who are most likely to benefit from the treatment. Efficacy can also be established according to validated endpoints within a timeframe that is compatible with further development of the therapy.

Challenges ahead and prospective developments

In HSC gene therapy, despite advances in vector design, determining the actual risk of insertional mutagenesis in clinical trials — and how it can be managed at the level of the individual patient or overcome by further improvements in vector design — remain major targets for future progress. Several experimental models have been developed to assess and rank the relative genotoxicity of different vector types and designs. Despite providing the rationale to advance new vectors to clinical testing, they fail to provide a quantitative prediction of actual oncogenic risk in the clinical setting. It is still difficult to establish the actual risk of genotoxicity for a given vector in a given trial, especially considering the low number of patients that have been treated with vectors to date, the longer follow-up periods that are required, and the possibility that the disease background could concur to increase the risk. In addition, it remains unclear whether molecular monitoring can predict the progression to malignancy of aberrant clones in individual patients. As long as the chosen vectors are able to integrate throughout the genome, subject to preferences dictated by the parental virus and particle composition, the risk of gene activation and inactivation at insertion sites can be mitigated but not abrogated. Although the improved vector types and designs that are currently being tested could alleviate some of the concerns that surround oncogene activation, which represents a high-risk dominant mutagenic event, the disruption of tumour-suppressor genes remains possible. These less frequent and recessive events might reveal their consequences only after a longer latency period or through testing in large-scale trials.

If efficacy and safety can be established in greater numbers of patients and over longer follow-up periods, HSC gene therapy might eventually challenge the dominance of allogeneic HSC transplantation as a first-line therapeutic option in monogenic diseases for which HSC transplantation is beneficial^{5,65}. In fact, the use of autologous HSCs makes the treatment

BOX 2

New adeno-associated virus vector serotypes

Major hurdles for liver-directed gene transfer include alleviating the inflammatory response to intravenous administration of high-dose viral particles and bypassing pre-existing immunity to viral components. Hepatocyte targeting must also be improved and long-term transgene expression should be established^{33,34}. In a recent clinical trial of liver-directed haemophilia B gene therapy, an adeno-associated virus (AAV) vector that had shown improved liver tropism (provided by the AAV8 serotype) and enhanced transduction efficiency (obtained by packaging a self-complementary genome) in animal models, resulted in the long-term expression of factor IX (FIX) in most patients⁴⁰ as long as a transient immunosuppressive corticosteroid treatment was promptly administered to those who were developing signs of hepatocellular injury. Whether the successful outcome can be attributed to the use of the new vector design remains to be established because the dose–response relationship and tropism of the different vector serotypes might differ between humans and animal models⁶⁴. Another ongoing trial of haemophilia B gene therapy is using the same AAV8 serotype to express a hyperactive FIX transgene. Preliminary reports indicate that sustained FIX expression occurs only in some patients, despite the administration of oral corticosteroids on signs of hepatocellular injury¹¹⁰. It is possible that the immunosuppressive treatment was administered too late to be effective at preventing the clearance of modified hepatocytes. Alternatively, certain aspects of vector design and manufacturing might be missed by current methods of characterization, which could influence the robustness of the therapy.

potentially available to all patients. HSC transplantation, meanwhile, is available only to patients for whom an HLA-matched or compatible donor can be found. HSC gene therapy can also substantially lower morbidity because it abrogates the risk of graft-versus-host disease and abolishes the need for immune suppression after treatment. Moreover, it could permit the use of milder preconditioning treatments, as partial chimaerism through the presence of gene-corrected cells might be sufficient to correct the disease and spare the patients the risk of myelo- and lymphoablative preconditioning treatments. In addition, genetic engineering of HSCs might facilitate new modes of treatment if the present outcome of HSC transplantation is unsatisfactory. This could involve increasing the therapeutic gene dosage to a level that is higher than the level that is provided by normal donor cells (as previously described for metachromatic leukodystrophy) or delivering transgene activity to selected tissues or sites of disease (such as the CNS in metachromatic leukodystrophy). HSCs and their progeny could also be equipped with exogenous genetic information to better fight cancer or chronic infection. As therapeutic options become available for otherwise incurable diseases, the genomes of newborn babies could be screened to allow early treatment of these conditions. To fulfil these predictions, *ex vivo* genetic modification and culturing of HSCs and their progenitors must become more robust and reliable. This should be coupled with the improved maintenance and expansion of the treated cells, which would support faster haematopoietic recovery in preconditioned patients and increase the clonal composition, resilience and stability of the engineered graft — thereby improving the short-term and long-term safety of the procedure. When therapeutic benefit is reached on establishing a partial chimaerism, less-toxic conditioning regimens based on biological agents could also be applied, which would spare patients both the acute and long-term toxicity of current chemotherapy-based regimens⁶⁶.

In T-cell immunotherapy for cancer, the number of CARs in clinical testing is growing quickly, although most of the trials focus on treating B-cell malignancies using CARs that target B-cell surface antigens. Notably, this approach depletes not just the disease-causing malignant clones but also almost all B cells in the patient. Although the lack of B cells can be remedied by the infusion of immunoglobulins, depletion of other cell lineages might not be as manageable. This issue could limit the use of other lineage-specific antigens as CAR targets. In addition, the clearance of large tumour masses observed in these trials was accompanied by an acute and often severe syndrome — even requiring intensive care — that followed the massive release of cytokines from on-target activated T cells.

Given the remarkable benefits that have been observed in some patients, further CARs are likely to be designed and tested^{42,55}. However, important questions still need to be addressed. For example, what CAR design will achieve the greatest efficacy with minimal toxicity? How can the toxicities deriving from the CAR-mediated recognition of healthy tissues that express target antigens at low levels be controlled in a timely and effective way? Can adoptive T-cell therapies be as effective in solid cancers as they seem to be in some lymphoid malignancies, or might features of the tumour stroma inhibit their activity? Although more challenging than transferring CAR genes, transferring TCR genes could be more effective in the long term. TCR gene transfer could also be a better fit for cancers that tend not to accumulate passenger mutations and might be less responsive to checkpoint blockade drugs. By targeting cancer-driver mutations, TCR gene transfer might be less prone to the induction of resistance than CAR-gene transfer, which has a more limited range of targets that must be surface antigens and are not always drivers of the disease. TCR gene editing^{67,68} can redirect specificity of a T cell towards a new antigen by disrupting endogenous TCR genes then introducing an exogenous TCR or CAR. It could therefore become a powerful strategy to avoid TCR mispairing, which occurs when the same cell expresses two different TCRs. Moreover, in driving the biological responses of the genetically modified T cells, it suppresses confounding endogenous TCR signal transduction. When combined with growing evidence to support the substantial efficacy of immune checkpoint blockade therapy, the outcome of recent CAR T-cell-gene transfer trials suggests that immunotherapy

could become a new pillar of cancer therapy that has the potential to eradicate the disease⁶⁹.

In some of the liver-directed, AAV-based, gene-therapy trials discussed in this Review, the consequences of a delayed cell-mediated immune response that targets transduced hepatocytes could be controlled. However, this remains an area of close scrutiny because the factors that might concur to trigger the response and regulate its timing are not fully understood³⁴. For instance, how does the occurrence of this delayed hepatocellular toxicity relate to the administered vector dose, potency and type? And is it related to possible concomitant triggers such as transient sub-clinical liver injury, or inflammation that recruits memory T cells at a time when viral antigen is still being presented by the transduced hepatocytes? Another unresolved aspect concerns the durability of transgene expression over an extended period of time. Although AAV-mediated FIX expression remains stable over the life of treated dogs, it is difficult to predict whether this also applies to humans, who have longer lifespans. A clear understanding of the molecular forms that underlie the long-term association of AAV DNA with the hepatocyte nucleus — and their relative contribution to transgene expression — will surely provide insight. In preclinical studies, a fraction of the AAV genome is reported to have integrated into the chromatin of hepatocytes, although it is unclear how much it could contribute to long-term gene expression³⁵. Although integration ensures long-term expression of genes, even in the face of cell replication, there are concerns over its safety^{70,71,72}. Approaches must be developed that enable patients with high concentrations of AAV-neutralizing antibodies to access AAV-mediated gene therapy and to allow re-administration of gene therapy, possibly by exploiting alternative capsid composition^{34,63,73,74}.

As confidence grows in the choice of vector type and dosing, patient selection criteria can be used to target therapies to individuals who are less likely to mount AAV-directed immune responses. New vector serotypes will ensure that AAV-based gene therapy becomes more widely applicable. For instance, it should be possible to tackle haemophilia A — the most common and challenging form of haemophilia — as well as several other metabolic and storage diseases that affect the liver and other peripheral organs, including muscle and the heart. The remarkable ability of AAV gene therapy to establish a long-term clinical benefit through a simple, well-tolerated intravenous infusion, combined with the ease of manufacturing a vector-based medicinal product (instead of a personalized cell product that is necessary for *ex vivo* gene therapy), could set the stage for its rapid commercial development and broad market distribution. AAV gene therapy might also be suitable for addressing unmet medical needs in countries with less well-developed health-care systems because a single intervention could alleviate the burden of providing and accessing lifelong replacement therapy and medical care. Other vector platforms, including lentiviral vectors, are also being developed for liver-directed gene therapy⁷⁵ and could eventually complement AAV-based delivery to broaden the access of patients to — and the diseases targeted by — this promising type of gene therapy.

Targeted gene editing

Another important boost to the gene therapy renaissance has emerged from advances in gene-targeting technologies⁷⁶. The ability to generate artificial DNA endonucleases that bind specifically to a DNA sequence of choice and induce a double-strand break (DSB) is making targeted genome editing more efficient and much easier to undertake, which brings the long-sought goals of somatic gene disruption, targeted transgene integration and *in situ* gene correction within the reach of gene therapy. All of these strategies entail a 'hit-and-run' mechanism that requires only transient expression of the nuclease complex and, in some cases, a repair template to modify the genome permanently. Because they target a selected region of the genome, these strategies abrogate the risk of insertional mutagenesis and poorly controlled transgene expression that is associated with conventional gene-replacement strategies.

Gene targeting has recently entered clinical testing through the adoptive T-cell therapy of patients infected with HIV⁷⁷. Artificial zinc-finger nucleases (ZFNs) were used to disrupt the gene that encodes CCR5, a cellular

co-receptor for HIV, in T cells grown *ex vivo* from each patient — with the aim of making these cells resistant to infection with the virus before their reinfusion. Gene disruption was achieved by creating a DNA DSB in an exon of the *CCR5* gene because repair of DSBs by the non-homologous end-joining (NHEJ) pathway leads to the loss or insertion of bases during end rejoining, which inactivates the coding sequence. The trial reported both proof of safety and long-term persistence of the engineered cells *in vivo*. It also provided an indication of efficacy through a trend for the positive selection of cells with disrupted *CCR5* alleles, as well as an observation in some patients of improved control of viral replication on the scheduled interruption of antiviral therapy⁷⁷. To achieve definitive viral clearance, however, the fraction and long-term maintenance of infused cells bearing biallelic *CCR5* disruption might need to be increased. This could be achieved by performing gene editing in self-renewing progenitor cells, such as HSCs⁷⁸ or T memory stem cells, because these cells offer better potential for achieving this level of enhanced and long-term reconstitution of the T-cell compartment, especially if combined with preconditioning to deplete endogenous non-modified cells. Targeted gene disruption can also be used to relieve the repression of an endogenous gene whose expression might compensate for the dysfunction of another gene. For instance, inactivation of the transcriptional repressor BCL11A in erythroid progenitor cells could reactivate fetal γ -globin expression to compensate for or counteract a dysfunctional β -globin chain in β -thalassaemia or sickle-cell disease, respectively⁷⁹. However, a higher efficiency of biallelic knockout is required to fulfil these goals when there is no mechanism by which to amplify selectively the genetically modified cells.

Gene editing becomes even more ambitious when aiming to replace a target sequence with an exogenous version of choice by exploiting the homology-directed repair (HDR) pathway of DNA DSBs. To achieve this, artificial nucleases and an exogenous DNA template bearing homology to the target site and comprising the new sequence⁷⁶ must be delivered to the cell. The approach has great potential for use in *ex vivo* gene therapy because the targeted integration of an expression cassette into a preselected genomic 'safe harbour'^{80,81} or the *in situ* reconstitution of a mutant gene would ensure robust and predictable expression — closely recapitulating the endogenous expression control in the latter — without the risk of insertional mutagenesis^{82–84}. Several hurdles must be overcome before these strategies can be fully exploited. This is because the efficiency of HDR-mediated genome editing remains low in most primary cell types of relevance to gene therapy, such as HSCs⁸². In addition, it is challenging to achieve the safe and feasible clinical translation of cell-therapy products when having to rely on selection and extensive *ex vivo* amplification of a few edited cell clones. The cellular response to DNA DSBs varies according to cell type and cell cycle and growth statuses, and ranges from repair by the different pathways to differentiation or apoptosis⁸⁵. Overall, how the cell chooses between NHEJ and HDR is poorly understood. Alternative mechanisms for HDR that might function outside the S/G2-phases of the cell cycle are also emerging.

Multiplexed, targeted genome editing is now easily achievable through the application of clustered regularly interspaced short palindromic repeat (CRISPR)/Cas9 RNA-based nucleases. These nucleases can be rapidly and easily adapted to seek out any DNA target site by designing an RNA guide instead of generating a protein-based sequence recognition motif for each target⁸⁶. Consequently, multiple applications have been found for targeted genome editing in experimental and preclinical models. Translating these applications to the clinic will, however, require thorough assessment of the off-target activity of the selected nuclease^{87,88} and optimization of the therapy. Gene addition by HDR or conventional vectors can also be combined with the targeted disruption of another cellular gene to augment effector function, such as by relieving inhibitory checkpoints when editing T-cell specificity⁸⁹.

Gene-editing strategies are also being developed for use in direct *in vivo* applications, including liver-directed gene therapy. Studies in mice have reported the reconstitution of FIX expression in haemophilia B mice by AAV-mediated delivery of ZFNs and a repair template⁹⁰, the targeted disruption of the cholesterol regulatory gene *Pcsk9* by AAV-mediated

BOX 3

Pricing gene therapy

The first gene therapy to be commercially approved in the Western world was alipogene tiparvovec (also known as Glybera). This muscle-directed adeno-associated virus 1-based gene therapy was granted marketing authorization in the European Union in 2012 for the treatment of a rare form of familial dyslipidaemia¹²¹. The market price was recently set at €1 million (US\$1.1 million) per treatment¹²². If forthcoming gene therapies are also sold at such a high price, they will challenge the standard reimbursement policies of governments and insurance companies. This high price reflects the cost of preclinical development, manufacturing and distribution of the new medicine, especially for *ex vivo* gene therapies, which are highly personalized and require individualized manufacturing. However, gene therapies have the potential to deliver a substantial, long-lasting benefit to the patient on a single administration, which may offset the cost of the standard treatment of the condition and its complications. Nonetheless, a single upfront payment model for gene therapy may not be sustainable. Approaches that spread the payment over several years should be considered, which could be linked to the successful outcome of the therapy. For example, a pay-for-performance strategy has been proposed¹²³ that makes the payment of instalments dependent on improved patient health, as determined by objective biomarkers. The risks are then shared between the health-care provider and insurer and the cost of treatment is more closely commensurate to the actual benefits delivered to the patient. In addition, the sustainability of gene therapies could be improved by adapting regulatory and manufacturing requirements to accommodate the unique features of these medicines and by facilitating their accessibility and distribution without detriment to their safety. A flexible platform-based approval and registration strategy should be considered, especially when developing gene and cell therapies that must be adapted to each individual, as is the case for the transfer of T-cell antigen receptor genes that target tumour neoantigens into autologous T cells of people with cancer^{44–47}.

delivery of CRISPR/Cas9 (ref. 91), and the repair of the *Fah* mutation in tyrosinaemia (an inability to break down the amino acid tyrosine) by hydrodynamic plasmid delivery of CRISPR/Cas9 and a template⁹². Major hurdles to the further development of *in vivo* gene editing include the safe and clinically suitable delivery of the editing machinery, which should act transiently without inducing cellular toxicity and immunogenicity. All artificial nucleases in current use employ one or more domains derived from prokaryotes — in some cases, common bacterial pathogens or saprophytes. Therefore, when delivered *in vivo*, there is a risk that such nucleases will generate or encounter a preexisting cell-mediated immunity, especially if they are expressed over the long term. More recently, an endonuclease-independent gene-targeting strategy was demonstrated in mice. A transgene was targeted to the albumin locus of hepatocytes by administering a promoterless AAV vector with homology to the highly transcribed albumin gene⁹³. A small proportion of the vector integrated at the albumin locus and became expressed from its endogenous promoter, which removed the requirement for both endonucleases and the transfer of a promoter within the vector. The efficiency and underlying mechanism of this and the other types of *in vivo* gene editing discussed in this Review are still to be fully determined.

Retinal gene therapy

The remarkable benefits that gene replacement can provide to patients with severe degenerative diseases were first highlighted by retinal gene therapy. Subretinal administration of AAV serotype 2 (AAV2)-mediated

gene therapy in patients with type 2 Leber congenital amaurosis (LCA), an inherited retinal dystrophy that causes loss of vision at an early age, led to improved visual acuity in several young patients in three independent trials^{94–96}. In two of these trials, the patients lost the benefit after a follow-up period of 2–3 years. However, sustained benefit was reported after a similar follow-up period in the third trial⁹⁷, which has now progressed to phase III testing. The reasons for the different outcome of these trials, all of which used an AAV2-derived capsid, remain unclear⁹⁸. The progressive degenerative nature of type 2 LCA poses a challenge for delivering extended therapeutic benefits because the small number of photoreceptors that are rescued by the therapy can eventually be overcome by non-cell-autonomous changes in the tissue. As discussed for liver-directed therapy, there might also be subtle differences in vector design and manufacturing that affect the extent of *in vivo* gene transfer, the inflammatory response at the site of delivery and the level of transgene reconstitution in transduced cells. The availability of AAV vectors with higher potencies, which would allow safer dose escalation and enhanced transduction, and more stringent tropism for the relevant targets, should help to overcome these limitations.

Other relevant developments

There are several other advances in the gene-therapy field that could not be discussed in detail in this Review. These include applications in neurodegenerative diseases that have reached the clinical-testing stage. For example, good safety but limited efficacy has been demonstrated for the delivery of transgenes to the brain by AAV or lentiviral vectors^{99–101}. Increasing the administered vector dose and optimizing the vehicle, cargo and study design are likely to lead to further advances.

Oncolytic viruses, which infect and kill cancer cells, have been in clinical use for some time. Although they can deliver robust and clinically relevant anticancer activity, these viruses are still being used as part of combination therapies¹⁰². Oncolytic therapies exploit viral replication and the induction of an immune response against infected cells — conditions that are normally offset in gene-therapy strategies. Intriguingly, their efficacy is likely to be augmented by adding a transgene cargo that improves the induction of anticancer immunity¹⁰³.

Adenoviral vectors of simian origin are also being assessed for their ability to induce humoral and cellular immunity through vaccination: encouraging results have already been seen in emerging or widespread infectious diseases that have long resisted conventional attempts^{104,105}. In preclinical models, the gene-based delivery of antibody therapy or prophylaxis is being explored to establish an *in vivo* stable and robust source of large quantity of antibodies with optimal specificity to treat or prevent infection¹⁰⁶.

Future outlook

Gene therapy could be poised to become an important new approach for the third millennium because its reach extends well beyond that of conventional drugs. Gene therapy enables the targeted delivery of information-rich gene-based cassettes that facilitate the stable, sustained and regulated expression of biological agents. Furthermore, when combined with cell therapy, it turns cells into smart vehicles for targeted gene delivery. As exemplified in the studies discussed in this Review, gene therapy directs powerful biological processes towards the goals of disease correction, tissue repair and regeneration. For instance, the stability, fidelity and amplification of the delivered therapeutic can be guaranteed by transferring information by genetic mechanisms. The homing and trafficking mechanisms of cells in the human body can be used to target gene-based therapeutics to specific tissues and disease sites. Gene therapy also makes use of the regenerative potential of stem cells and transplantation as well as the biological weapon of immunity, which is exploited for the specific elimination of transformed or infected cells. By taking advantage of these inbuilt biological capabilities, gene therapy has the potential to address the substantial unmet medical needs of both rare and common severe diseases, which will benefit both patients and — more broadly — society. Major challenges must still be addressed before this promise can be

realized. For example, the efficacy and safety of gene-transfer vectors should be improved by further engineering their design and composition, which could include combining the biological features of different viruses with synthetic molecules. These advances will enable vectors to target tissues and cell types precisely^{64,73,107}, and overcome cellular restrictions on gene transduction and bypass sensors of exogenous nucleic acids. They will also help vectors to avoid activating the innate and adaptive immune system. Overall, the changes will also ensure that transgene expression is reproducible, robust, occurs over an extended period and closely mimics the endogenous pattern of expression (when gene replacement is performed). Improvements in vector manufacturing and characterization will allow the standardization and comparative assessment of vector performance between trials. The rate and specificity of *in situ* gene correction and editing, the integration of vectors at safe genomic harbours, and allele-specific silencing by artificial nucleases and epigenetic modifiers represent further opportunities for improving gene-therapy strategies. Deeper understanding of disease pathogenesis in inherited, multigenic or acquired conditions will enable the development of new gene-based treatment strategies. Because gene transfers employ ‘live’ biological drugs of unprecedented complexity that have the potential to induce extended effects on patients and their germ lines, long-term surveillance and precautionary measures must be taken while their use is being pioneered. Moreover, current — limited — understanding of the regulation of stem cells, tissue regeneration and immune-response checkpoints constrains the capacity for intervention and raises concerns about the untoward effects of manipulation.

From a clinical standpoint, the bedside delivery of gene and cell therapies calls for multidisciplinary expertise and, in some cases, advanced cell processing at the clinical-treatment site. Biological readouts must also be developed to monitor the safety and efficacy of therapies. As the first gene therapies progress from registration to marketing, both the pharmaceutical sector and regulatory agencies are being engaged to help define appropriate quality standards for manufacturing and release and to build suitable pipelines for supplying such highly personalized therapies. From a societal standpoint, the complexity and cost of manufacturing and supplying ‘live’ biological drugs in conventional health-care systems will challenge the sustainability of these therapies and require creative cost-reimbursement policies that enable all patients to benefit from them (Box 3).

Finally, from an ethics standpoint, it is important to consider whether medicine should surrender to the rule of technology or commit to a more responsible steering of the course of progress. For instance, the avenues that are being opened to intervention by emerging technologies could undermine our self-perception and self-determination as we end up viewing ourselves as the evolutionary product of DNA that has become self-conscious and can edit itself to shape its progeny as and when it desires. The call for a moratorium on applying genome editing to human germline cells highlights forthcoming ethical dilemmas^{108,109}. Science might turn again to the ancient roots of Western culture to learn from its wisdom, such as the “Know thyself” inscription on the Temple of Apollo at Delphi in ancient Greece — a warning to recognize our limits. Modern philosophy has trained us to exercise criticism when assessing the truth and certainty of knowledge, which is limited a priori and enables learning relationships and making predictions but not uncovering the nature of things. Those predictions allow us to develop strategies that can alleviate human suffering from disease, thereby providing our endeavours with a translational framework that can guide our choices and justify them from an ethics standpoint. ■

Received 15 April; accepted 24 August 2015.

1. Naldini, L. *Ex vivo* gene transfer and correction for cell-based therapies. *Nature Rev. Genet.* **12**, 301–315 (2011).
2. Hacein-Bey-Abina, S. *et al.* Efficacy of gene therapy for X-linked severe combined immunodeficiency. *N. Engl. J. Med.* **363**, 355–364 (2010).
3. Aiuti, A. *et al.* Gene therapy for immunodeficiency due to adenosine deaminase deficiency. *N. Engl. J. Med.* **360**, 447–458 (2009).
4. Ferrua, F., Brigida, I. & Aiuti, A. Update on gene therapy for adenosine deaminase-deficient severe combined immunodeficiency. *Curr. Opin. Allergy Clin. Immunol.* **10**, 551–556 (2010).

5. Fischer, A., Hacein-Bey-Abina, S. & Cavazzana-Calvo, M. 20 years of gene therapy for SCID. *Nature Immunol.* **11**, 457–460 (2010).
A comprehensive review of the therapeutic potential, risks and limitations of HSC-based SCID gene therapy using γ -RV by some of its pioneers; see also refs 3 and 7.
6. Gaspar, H. B. *et al.* Hematopoietic stem cell gene therapy for adenosine deaminase-deficient severe combined immunodeficiency leads to long-term immunological recovery and metabolic correction. *Sci. Transl. Med.* **3**, 97ra80 (2011); erratum **5**, 168er1 (2013).
7. Gaspar, H. B. *et al.* Long-term persistence of a polyclonal T cell repertoire after gene therapy for X-linked severe combined immunodeficiency. *Sci. Transl. Med.* **3**, 97ra79 (2011).
8. Boztug, K. *et al.* Stem-cell gene therapy for the Wiskott–Aldrich syndrome. *N. Engl. J. Med.* **363**, 1918–1927 (2010).
9. Candotti, F. *et al.* Gene therapy for adenosine deaminase-deficient severe combined immune deficiency: clinical comparison of retroviral vectors and treatment plans. *Blood* **120**, 3635–3646 (2012).
10. Kang, E. M. *et al.* Retrovirus gene therapy for X-linked chronic granulomatous disease can achieve stable long-term correction of oxidase activity in peripheral blood neutrophils. *Blood* **115**, 783–791 (2010).
11. Hacein-Bey-Abina, S. *et al.* Insertional oncogenesis in 4 patients after retrovirus-mediated gene therapy of SCID-X1. *J. Clin. Invest.* **118**, 3132–3142 (2008).
12. Howe, S. J. *et al.* Insertional mutagenesis combined with acquired somatic mutations causes leukemogenesis following gene therapy of SCID-X1 patients. *J. Clin. Invest.* **118**, 3143–3150 (2008).
13. Stein, S. *et al.* Genomic instability and myelodysplasia with monosomy 7 consequent to EVI1 activation after gene therapy for chronic granulomatous disease. *Nature Med.* **16**, 198–204 (2010).
14. Braun, C. J. *et al.* Gene therapy for Wiskott–Aldrich syndrome—long-term efficacy and genotoxicity. *Sci. Transl. Med.* **6**, 227ra233 (2014).
15. Kang, H. J. *et al.* Retroviral gene therapy for X-linked chronic granulomatous disease: results from phase I/II trial. *Mol. Ther.* **19**, 2092–2101 (2011).
16. Aiuti, A. *et al.* Lentiviral hematopoietic stem cell gene therapy in patients with Wiskott–Aldrich syndrome. *Science* **341**, 1233151 (2013).
In this study, vector insertional analyses in patients show data consistent with improved safety of lentiviral vectors versus γ -RVs while achieving similarly effective disease correction; see also ref. 17.
17. Hacein-Bey-Abina, S. *et al.* Outcomes following gene therapy in patients with severe Wiskott–Aldrich syndrome. *J. Am. Med. Assoc.* **313**, 1550–1563 (2015).
18. Hacein-Bey-Abina, S. *et al.* A modified γ -retrovirus vector for X-linked severe combined immunodeficiency. *N. Engl. J. Med.* **371**, 1407–1417 (2014).
19. Cavazzana-Calvo, M. *et al.* Transfusion independence and *HMGA2* activation after gene therapy of human β -thalassaemia. *Nature* **467**, 318–322 (2010).
20. Cartier, N. *et al.* Hematopoietic stem cell gene therapy with a lentiviral vector in X-linked adrenoleukodystrophy. *Science* **326**, 818–823 (2009).
The first trial of HSC gene therapy performed with lentiviral vectors shows data consistent with stable HSC transduction, with long-term safety and efficacy revealed in the follow-up paper (see ref. 21).
21. Cartier, N. *et al.* Lentiviral hematopoietic cell gene therapy for X-linked adrenoleukodystrophy. *Methods Enzymol.* **507**, 187–198 (2012).
22. Biffi, A. *et al.* Lentiviral hematopoietic stem cell gene therapy benefits metachromatic leukodystrophy. *Science* **341**, 1233158 (2013).
This study highlights the potential of genetic engineering by achieving the stable reconstitution of haematopoiesis in which up to 90% of cells are gene corrected and overexpress the transgene, which provides therapeutic benefit when conventional HSC transplantation is less satisfactory.
23. Notarangelo, L. D. *et al.* Primary immunodeficiencies: 2009 update. *J. Allergy Clin. Immunol.* **124**, 1161–1178 (2009).
24. Kemp, S., Berger, J. & Aubourg, P. X-linked adrenoleukodystrophy: clinical, metabolic, genetic and pathophysiological aspects. *Biochim. Biophys. Acta* **1822**, 1465–1474 (2012).
25. Gieselmann, V. & Krageloh-Mann, I. Metachromatic leukodystrophy—an update. *Neuropediatrics* **41**, 1–6 (2010).
26. Gennery, A. R. *et al.* Transplantation of hematopoietic stem cells and long-term survival for primary immunodeficiencies in Europe: entering a new century, do we do better? *J. Allergy Clin. Immunol.* **126**, 602–610 (2010).
27. Krageloh-Mann, I. *et al.* Juvenile metachromatic leukodystrophy 10 years post transplant compared with a non-transplanted cohort. *Bone Marrow Transplant.* **48**, 369–375 (2013).
28. Copelan, E. A. Hematopoietic stem-cell transplantation. *N. Engl. J. Med.* **354**, 1813–1826 (2006).
29. Biffi, A. *et al.* Correction of metachromatic leukodystrophy in the mouse model by transplantation of genetically modified hematopoietic stem cells. *J. Clin. Invest.* **113**, 1118–1129 (2004).
30. Biffi, A. *et al.* Gene therapy of metachromatic leukodystrophy reverses neurological damage and deficits in mice. *J. Clin. Invest.* **116**, 3070–3082 (2006).
31. Capotondo, A. *et al.* Brain conditioning is instrumental for successful microglia reconstitution following hematopoietic stem cell transplantation. *Proc. Natl Acad. Sci. USA* **109**, 15018–15023 (2012).
32. Mingozzi, F. & High, K. A. Therapeutic *in vivo* gene transfer for genetic disease using AAV: progress and challenges. *Nature Rev. Genet.* **12**, 341–355 (2011).
33. Nayak, S. & Herzog, R. W. Progress and prospects: immune responses to viral vectors. *Gene Ther.* **17**, 295–304 (2010).
34. Mingozzi, F. & High, K. A. Immune responses to AAV vectors: overcoming barriers to successful gene therapy. *Blood* **122**, 23–36 (2013).
35. Grieger, J. C. & Samulski, R. J. Adeno-associated virus vectorology, manufacturing and clinical applications. *Methods Enzymol.* **507**, 229–254 (2012).
36. High, K. H., Nathwani, A., Spencer, T. & Lillicrap, D. Current status of haemophilia gene therapy. *Haemophilia* **20** (suppl. 4), 43–49 (2014).
37. Bernthorpe, E. & Shapiro, A. D. Modern haemophilia care. *Lancet* **379**, 1447–1456 (2012).
38. Manno, C. S. *et al.* Successful transduction of liver in hemophilia by AAV-Factor IX and limitations imposed by the host immune response. *Nature Med.* **12**, 342–347 (2006); erratum **12**, 592 (2006).
The first clinical data to show the safety and potential efficacy of liver-directed AAV gene transfer, which was unexpectedly abrogated by an immune response against viral capsids (as detailed in ref. 39).
39. Mingozzi, F. *et al.* CD8⁺ T-cell responses to adeno-associated virus capsid in humans. *Nature Med.* **13**, 419–422 (2007).
40. Nathwani, A. C. *et al.* Long-term safety and efficacy of factor IX gene therapy in hemophilia B. *N. Engl. J. Med.* **371**, 1994–2004 (2014).
This AAV8-based trial was first to report stable FIX expression at therapeutic levels and also first to overcome the detrimental effect of the immune response to viral capsids by corticosteroid administration.
41. Rosenberg, S. A. & Restifo, N. P. Adoptive cell transfer as personalized immunotherapy for human cancer. *Science* **348**, 62–68 (2015).
A comprehensive review of the clinical development and potentially transformative impact of adoptive T-cell therapy on cancer by one of its pioneers; see also ref. 42.
42. Maus, M. V. *et al.* Adoptive immunotherapy for cancer or viruses. *Annu. Rev. Immunol.* **32**, 189–225 (2014).
43. Dudley, M. E. *et al.* Cancer regression and autoimmunity in patients after clonal repopulation with antitumor lymphocytes. *Science* **298**, 850–854 (2002).
44. Gubin, M. M. *et al.* Checkpoint blockade cancer immunotherapy targets tumour-specific mutant antigens. *Nature* **515**, 577–581 (2014).
45. Yadav, M. *et al.* Predicting immunogenic tumour mutations by combining mass spectrometry and exome sequencing. *Nature* **515**, 572–576 (2014).
46. Schumacher, T. N. & Schreiber, R. D. Neoantigens in cancer immunotherapy. *Science* **348**, 69–74 (2015).
A timely review on the origin and nature of tumour neoantigens and how they can be identified and potentially exploited for targeted T-cell gene therapy in the clinical setting.
47. Rizvi, N. A. *et al.* Cancer immunology. Mutational landscape determines sensitivity to PD-1 blockade in non-small cell lung cancer. *Science* **348**, 124–128 (2015).
48. Hunder, N. N. *et al.* Treatment of metastatic melanoma with autologous CD4⁺ T cells against NY-ESO-1. *N. Engl. J. Med.* **358**, 2698–2703 (2008).
49. Johnson, L. A. *et al.* Gene therapy with human and mouse T-cell receptors mediates cancer regression and targets normal tissues expressing cognate antigen. *Blood* **114**, 535–546 (2009).
50. Robbins, P. F. *et al.* A pilot trial using lymphocytes genetically engineered with an NY-ESO-1-reactive T-cell receptor: long-term follow-up and correlates with response. *Clin. Cancer Res.* **21**, 1019–1027 (2015).
51. Kochenderfer, J. N. *et al.* Chemotherapy-refractory diffuse large B-cell lymphoma and indolent B-cell malignancies can be effectively treated with autologous T cells expressing an anti-CD19 chimeric antigen receptor. *J. Clin. Oncol.* **33**, 540–549 (2015).
52. Brentjens, R. J. *et al.* CD19-targeted T cells rapidly induce molecular remissions in adults with chemotherapy-refractory acute lymphoblastic leukemia. *Sci. Transl. Med.* **5**, 177ra38 (2013).
53. Maude, S. L. *et al.* Chimeric antigen receptor T cells for sustained remissions in leukemia. *N. Engl. J. Med.* **371**, 1507–1517 (2014).
54. Lee, D. W. *et al.* T cells expressing CD19 chimeric antigen receptors for acute lymphoblastic leukaemia in children and young adults: a phase 1 dose-escalation trial. *Lancet* **385**, 517–528 (2015).
55. June, C. H., Riddell, S. R. & Schumacher, T. N. Adoptive cellular therapy: a race to the finish line. *Sci. Transl. Med.* **7**, 280ps7 (2015).
56. Biffi, A. *et al.* Lentiviral vector common integration sites in preclinical models and a clinical trial reflect a benign integration bias and not oncogenic selection. *Blood* **117**, 5332–5339 (2011).
57. Deichmann, A. *et al.* Insertion sites in engrafted cells cluster within a limited repertoire of genomic areas after gammaretroviral vector gene therapy. *Mol. Ther.* **19**, 2031–2039 (2011).
58. Doulatov, S., Notta, F., Laurenti, E. & Dick, J. E. Hematopoiesis: a human perspective. *Cell Stem Cell* **10**, 120–136 (2012).
59. Gattinoni, L. Memory T cells officially join the stem cell club. *Immunity* **41**, 7–9 (2014).
60. Cieri, N. *et al.* IL-7 and IL-15 instruct the generation of human memory stem T cells from naive precursors. *Blood* **121**, 573–584 (2013).
61. Biasco, L. *et al.* *In vivo* tracking of T cells in humans unveils decade-long survival and activity of genetically modified T memory stem cells. *Sci. Transl. Med.* **7**, 273ra13 (2015).
62. Asokan, A., Schaffer, D. V. & Samulski, R. J. The AAV vector toolkit: poised at the clinical crossroads. *Mol. Ther.* **20**, 699–708 (2012).
63. Mingozzi, F. *et al.* Overcoming preexisting humoral immunity to AAV using capsid decoys. *Sci. Transl. Med.* **5**, 194ra92 (2013).
64. Lisowski, L. *et al.* Selection and evaluation of clinically relevant AAV variants in a xenograft liver model. *Nature* **506**, 382–386 (2014).
65. Kohn, D. B. Gene therapy outpaces haplo for SCID-X1. *Blood* **125**, 3521–3522 (2015).
66. Logan, A. C., Weissman, I. L. & Shizuru, J. A. The road to purified hematopoietic stem cell transplants is paved with antibodies. *Curr. Opin. Immunol.* **24**, 640–648 (2012).

67. Provasi, E. *et al.* Editing T cell specificity towards leukemia by zinc finger nucleases and lentiviral gene transfer. *Nature Med.* **18**, 807–815 (2012).
68. Torikai, H. *et al.* A foundation for universal T-cell based immunotherapy: T cells engineered to express a CD19-specific chimeric-antigen-receptor and eliminate expression of endogenous TCR. *Blood* **119**, 5697–5705 (2012).
69. Sharma, P. & Allison, J. P. The future of immune checkpoint therapy. *Science* **348**, 56–61 (2015).
70. Li, H. *et al.* Assessing the potential for AAV vector genotoxicity in a murine model. *Blood* **117**, 3311–3319 (2011).
71. Chandler, R. J. *et al.* Vector design influences hepatic genotoxicity after adeno-associated virus gene therapy. *J. Clin. Invest.* **125**, 870–880 (2015).
72. Nault, J.-C. *et al.* Recurrent AAV2-related insertional mutagenesis in human hepatocellular carcinomas. *Nature Genet.* <http://dx.doi.org/10.1038/ng.3389> (2015).
73. Martino, A. T. *et al.* Engineered AAV vector minimizes *in vivo* targeting of transduced hepatocytes by capsid-specific CD8⁺ T cells. *Blood* **121**, 2224–2233 (2013).
74. Kotterman, M. A. & Schaffer, D. V. Engineering adeno-associated viruses for clinical gene therapy. *Nature Rev. Genet.* **15**, 445–451 (2014).
75. Cantore, A. *et al.* Liver-directed lentiviral gene therapy in a dog model of hemophilia B. *Sci. Transl. Med.* **7**, 277ra28 (2015).
76. Cox, D. B., Platt, R. J. & Zhang, F. Therapeutic genome editing: prospects and challenges. *Nature Med.* **21**, 121–131 (2015).
77. Tebas, P. *et al.* Gene editing of CCR5 in autologous CD4 T cells of persons infected with HIV. *N. Engl. J. Med.* **370**, 901–910 (2014).
The first clinical testing of targeted gene disruption that showed the safety, persistence and survival advantage of T cells that have been genetically edited for resistance to HIV-1.
78. Li, L. *et al.* Genomic editing of the HIV-1 coreceptor CCR5 in adult hematopoietic stem and progenitor cells using zinc finger nucleases. *Mol. Ther.* **21**, 1259–1269 (2013).
79. Bauer, D. E. *et al.* An erythroid enhancer of BCL11A subject to genetic variation determines fetal hemoglobin level. *Science* **342**, 253–257 (2013).
80. Lombardo, A. *et al.* Site-specific integration and tailoring of cassette design for sustainable gene transfer. *Nature Methods* **8**, 861–869 (2011).
81. Rio, P. *et al.* Targeted gene therapy and cell reprogramming in Fanconi anemia. *EMBO Mol. Med.* **6**, 835–848 (2014).
82. Genovese, P. *et al.* Targeted genome editing in human repopulating haematopoietic stem cells. *Nature* **510**, 235–240 (2014).
This paper demonstrates differential permissiveness to targeted genome editing in haematopoietic stem and progenitor cells and provides a proof of concept for the *in situ* correction of SCID-X1 mutations in HSCs.
83. Hoban, M. D. *et al.* Correction of the sickle cell disease mutation in human hematopoietic stem/progenitor cells. *Blood* **125**, 2597–2604 (2015).
84. Osborn, M. J. *et al.* Fanconi anemia gene editing by the CRISPR/Cas9 system. *Hum. Gene Ther.* **26**, 114–126 (2015).
85. Jasin, M. & Rothstein, R. Repair of strand breaks by homologous recombination. *Cold Spring Harb. Perspect. Biol.* **5**, a012740 (2013).
86. Doudna, J. A. & Charpentier, E. Genome editing. The new frontier of genome engineering with CRISPR-Cas9. *Science* **346**, 1258096 (2014).
In this review, the researchers who pioneered the application of RNA-guided nucleases to genome engineering show how this transformative technique can make targeted genome editing easy.
87. Gabriel, R. *et al.* An unbiased genome-wide analysis of zinc-finger nuclease specificity. *Nature Biotechnol.* **29**, 816–823 (2011).
88. Tsai, S. Q. *et al.* GUIDE-seq enables genome-wide profiling of off-target cleavage by CRISPR-Cas nucleases. *Nature Biotechnol.* **33**, 187–197 (2015).
89. Beane, J. D. *et al.* Clinical scale zinc finger nuclease-mediated gene editing of PD-1 in tumor infiltrating lymphocytes for the treatment of metastatic melanoma. *Mol. Ther.* **23**, 1380–1390 (2015).
90. Li, H. *et al.* *In vivo* genome editing restores haemostasis in a mouse model of haemophilia. *Nature* **475**, 217–221 (2011).
The first study to show the feasibility of targeted genome editing *in vivo* by the AAV-mediated delivery of artificial nucleases and template.
91. Ran, F. A. *et al.* *In vivo* genome editing using *Staphylococcus aureus* Cas9. *Nature* **520**, 186–191 (2015).
92. Yin, H. *et al.* Genome editing with Cas9 in adult mice corrects a disease mutation and phenotype. *Nature Biotechnol.* **32**, 551–553 (2014).
93. Barzel, A. *et al.* Promoterless gene targeting without nucleases ameliorates haemophilia B in mice. *Nature* **517**, 360–364 (2015).
94. Simonelli, F. *et al.* Gene therapy for Leber's congenital amaurosis is safe and effective through 1.5 years after vector administration. *Mol. Ther.* **18**, 643–650 (2010).
95. Jacobson, S. G. *et al.* Improvement and decline in vision with gene therapy in childhood blindness. *N. Engl. J. Med.* **372**, 1920–1926 (2015).
96. Bainbridge, J. W. *et al.* Long-term effect of gene therapy on Leber's congenital amaurosis. *N. Engl. J. Med.* **372**, 1887–1897 (2015).
97. Testa, F. *et al.* Three-year follow-up after unilateral subretinal delivery of adeno-associated virus in patients with Leber congenital amaurosis type 2. *Ophthalmology* **120**, 1283–1291 (2013).
98. Wright, A. F. Long-term effects of retinal gene therapy in childhood blindness. *N. Engl. J. Med.* **372**, 1954–1955 (2015).
99. Leone, P. *et al.* Long-term follow-up after gene therapy for canavan disease. *Sci. Transl. Med.* **4**, 165ra163 (2012).
100. Tardieu, M. *et al.* Intracerebral administration of adeno-associated viral vector serotype rh.10 carrying human SGSH and SUMF1 cDNAs in children with mucopolysaccharidosis type IIIA disease: results of a phase I/II trial. *Hum. Gene Ther.* **25**, 506–516 (2014).
101. Palfi, S. *et al.* Long-term safety and tolerability of ProSavin, a lentiviral vector-based gene therapy for Parkinson's disease: a dose escalation, open-label, phase 1/2 trial. *Lancet* **383**, 1138–1146 (2014).
102. Miest, T. S. & Cattaneo, R. New viruses for cancer therapy: meeting clinical needs. *Nature Rev. Microbiol.* **12**, 23–34 (2014).
103. Lichty, B. D., Breitbach, C. J., Stojdl, D. F. & Bell, J. C. Going viral with cancer immunotherapy. *Nature Rev. Cancer* **14**, 559–567 (2014).
104. Ogwang, C. *et al.* Prime-boost vaccination with chimpanzee adenovirus and modified vaccinia Ankara encoding TRAP provides partial protection against *Plasmodium falciparum* infection in Kenyan adults. *Sci. Transl. Med.* **7**, 286re5 (2015).
105. Rampling, T. *et al.* A monovalent chimpanzee adenovirus ebola vaccine — preliminary report. *N. Engl. J. Med.* <http://dx.doi.org/10.1056/NEJMoa1411627> (2015).
106. Balazs, A. B. *et al.* Vectored immunoprophylaxis protects humanized mice from mucosal HIV transmission. *Nature Med.* **20**, 296–300 (2014).
107. Girard-Gagnepain, A. *et al.* Baboon envelope pseudotyped LVs outperform VSV-G-LVs for gene transfer into early-cytokine-stimulated and resting HSCs. *Blood* **124**, 1221–1231 (2014).
108. Baltimore, D. *et al.* Biotechnology. A prudent path forward for genomic engineering and germline gene modification. *Science* **348**, 36–38 (2015).
109. Bosley, K. S. *et al.* CRISPR germline engineering—the community speaks. *Nature Biotechnol.* **33**, 478–486 (2015).
110. Baxter. Baxalta reports continued progress on phase 1/2 clinical trial of BAX335, investigational gene therapy treatment for hemophilia B. *Baxter* http://www.baxter.com/news-media/newsroom/press-releases/2015/06_24_15_bax335.page (2015).
111. Montini, E. *et al.* Hematopoietic stem cell gene transfer in a tumor-prone mouse model uncovers low genotoxicity of lentiviral vector integration. *Nature Biotechnol.* **24**, 687–696 (2006).
112. Modlich, U. *et al.* Cell-culture assays reveal the importance of retroviral vector design for insertional genotoxicity. *Blood* **108**, 2545–2553 (2006).
113. Zychlinski, D. *et al.* Physiological promoters reduce the genotoxic risk of integrating gene vectors. *Mol. Ther.* **16**, 718–725 (2008).
114. Montini, E. *et al.* The genotoxic potential of retroviral vectors is strongly modulated by vector design and integration site selection in a mouse model of HSC gene therapy. *J. Clin. Invest.* **119**, 964–975 (2009).
This preclinical study highlights important features of vector design that affect genotoxicity and reveals strategies to alleviate it; the study was instrumental in promoting the clinical testing of improved vectors (see refs 112–113 for an *in vitro* assay that provides complementary information).
115. Modlich, U. *et al.* Insertional transformation of hematopoietic cells by self-inactivating lentiviral and gammaretroviral vectors. *Mol. Ther.* **17**, 1919–1928 (2009).
116. Zhou, S. *et al.* A self-inactivating lentiviral vector for SCID-X1 gene therapy that does not activate LMO2 expression in human T cells. *Blood* **116**, 900–908 (2010).
117. Zhou, S. *et al.* Mouse transplant models for evaluating the oncogenic risk of a self-inactivating XSCID lentiviral vector. *PLoS ONE* **8**, e62333 (2013).
118. Baum, C., Modlich, U., Gohring, G. & Schlegelberger, B. Concise review: managing genotoxicity in the therapeutic modification of stem cells. *Stem Cells* **29**, 1479–1484 (2011).
119. Amendola, M., Venneri, M. A., Biffi, A., Vigna, E. & Naldini, L. Coordinate dual-gene transgenesis by lentiviral vectors carrying synthetic bidirectional promoters. *Nature Biotechnol.* **23**, 108–116 (2005).
120. Greco, R. *et al.* Improving the safety of cell therapy with the TK-suicide gene. *Front. Pharmacol.* **6**, 95 (2015).
121. Melchiorri, D. *et al.* Regulatory evaluation of Glybera in Europe — two committees, one mission. *Nature Rev. Drug Discov.* **12**, 719 (2013).
122. Morrison, C. \$1-million price tag set for Glybera gene therapy. *Nature Biotechnol.* **33**, 217–218 (2015).
123. Brennan, T. A. & Wilson, J. M. The special case of gene therapy pricing. *Nature Biotechnol.* **32**, 874–876 (2014).

Acknowledgements L.N. apologizes to the many scientists whose contributions to the field could not be acknowledged owing to space limitations. He thanks past and present members of his laboratory and colleagues at the San Raffaele Telethon Institute for Gene Therapy (TIGET) and the San Raffaele Scientific Institute. L.N. is also grateful to the Telethon Foundation, the European Union (FP7 and ERC), the Italian Association for Cancer Research, and the Italian ministries of health and of scientific research for supporting his research.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The author declares competing financial interests: see go.nature.com/pnd2xw. Readers are welcome to comment on the online version of this paper at go.nature.com/pnd2xw. Correspondence should be addressed to L.N. (naldini.luigi@hsr.it).

Patient-centric trials for therapeutic development in precision oncology

Andrew V. Biankin^{1,2,3,4}, Steven Piantadosi⁵ & Simon J. Hollingsworth⁶

An enhanced understanding of the molecular pathology of disease gained from genomic studies is facilitating the development of treatments that target discrete molecular subclasses of tumours. Considerable associated challenges include how to advance and implement targeted drug-development strategies. Precision medicine centres on delivering the most appropriate therapy to a patient on the basis of clinical and molecular features of their disease. The development of therapeutic agents that target molecular mechanisms is driving innovation in clinical-trial strategies. Although progress has been made, modifications to existing core paradigms in oncology drug development will be required to realize fully the promise of precision medicine.

Insights into the molecular pathology of disease are creating opportunities for the development of therapies with durable clinical benefit while challenging the existing model of therapeutic development and clinical care^{1–3}. Large international consortia — such as the International Cancer Genome Consortium^{4,5} — are mapping the genomes of thousands of cancers to identify opportunities for prevention, early detection and treatment⁶. Although genomics is leading the way, high-throughput proteomics and metabolomics are following closely behind⁷. Such methodological advances have ushered in a new era of therapeutics that target specific molecular processes. Although there have been some dramatic successes^{8–17}, the overall strategy remains in its infancy¹⁸. The central premise of precision medicine is that matching a drug and its mechanism of action using a marker to select patients — a process often referred to as matching the right drug to the right patient — can offer greater potential for durable clinical benefits.

Initially, these targeted therapeutic agents followed the same clinical development pathway as cytotoxic chemotherapy, that is, based on tumour location and histopathology, driven by the notion that molecular aberrations were tumour specific. Efforts to advance this approach stalled because of the lack of efficacy data in patients with different cancer types that shared a molecular aberration, coupled with early observations that the functional importance of some aberrations varied between tumour types. However, the emergence of programmes that identified molecular targets and matched treatments to molecular subtypes — or segments — led to several reports^{19,20,21} that directly linked this approach to improvements in clinical outcome, irrespective of the organ in which the tumour originated. Although many were based on retrospective analyses of tumour samples, and not all reports were equally convincing²², the utility of broad molecular profiling to guide patients towards specific targeted therapies was established. Researchers moved quickly to implement this new paradigm. To meet emerging requirements, and enticed by the promise of clinical benefit, clinicians recognized that the established pathways of therapeutic development would need to change. However, the practical implications of implementing these changes in the clinic were unclear.

The drivers of precision medicine have been established and discussed elsewhere^{18,23,24}. However, fresh challenges for therapeutic

development are many and substantial. Fundamentally, a candidate treatment requires a strong platform of evidence to support its clinical testing and must be coupled with robust methods to identify appropriate patients (using molecular assays²⁵). Our appreciation of the molecular diversity of cancer and the ever-increasing number of molecular subtypes creates considerable complexity for the development of targeted drugs. When tested in trials of unselected participants, most targeted therapies reveal efficacy only if both the incidence of a responsive subpopulation and the effect size within the group is sufficiently high. Increasing the size of clinical trials to overcome this lack of enrichment yields minimal overall benefits at a cost that makes them unattractive and unaffordable to the community. Designing trials that feasibly evaluate both patient selection and drug efficacy is crucial, and it is essential to define the correct metrics to assess efficacy, particularly when the study needs to be small.

Principles and evolution of clinical trials

Clinical trials are most useful when they assess a potential therapeutic effect that is about the same size or slightly smaller than the effect of the natural variation that exists between individuals. When the variation between individuals enrolled in a trial influences a treatment only randomly, it can be ignored in a biological sense and controlled by replication. These dual strategies for controlling for variation embody the empirical and theoretical aspects of trials. For much of the history of clinical trials, the treatments under investigation were assumed to apply to anyone with the relevant clinically defined condition. Essentially, our understanding of biology suggested that treatments worked through common mechanisms that were set apart from random variation. This assumption was substantially correct for approaches such as cytotoxic chemotherapy that target generic disease mechanisms, and it enabled considerable progress to be made in treating cancer. Towards the end of the twentieth century, concerns arose regarding the potential inhomogeneity of therapeutic effects because of socio-political characteristics such as race or sex. Many clinical trials were designed and analysed to examine such differences. Although motivated by politics and social justice rather than scientific fact, only minimal changes were actually made to the design of such trials — which was probably appropriate given the

¹Wolfson Wohl Cancer Research Centre, Institute of Cancer Sciences, University of Glasgow, Glasgow, Scotland G61 1BD, UK. ²The Kinghorn Cancer Centre, Cancer Division, Garvan Institute of Medical Research, Sydney, New South Wales 2010, Australia. ³Department of Surgery, Bankstown Hospital, Sydney, New South Wales 2200, Australia. ⁴South Western Sydney Clinical School, Faculty of Medicine, University of New South Wales, Liverpool, New South Wales 2170, Australia. ⁵Samuel Oschin Comprehensive Cancer Institute, Cedars-Sinai Medical Center, Los Angeles, California 90095, USA. ⁶Innovative Medicines & Early Development Oncology, AstraZeneca, Cambridge Science Park, Cambridge CB4 0FZ, UK.

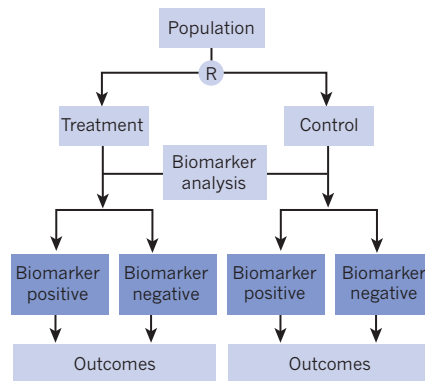
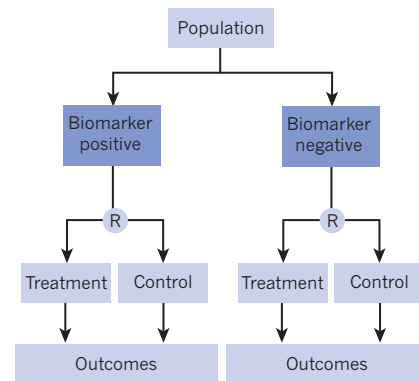
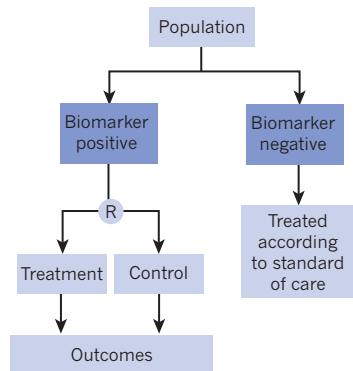
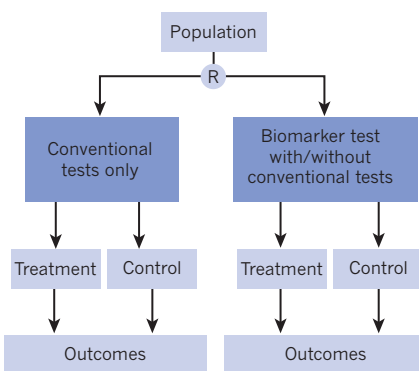
a Biomarker analysis within existing RCT**b Non-targeted RCT (stratified by biomarker)****c Targeted RCT****d Classical RCT**

Figure 1 | Randomized controlled trial designs for defining and testing precision-medicine strategies. **a**, Biomarker discovery is performed in a trial that is used to address a therapeutic question but patient recruitment and treatment allocation are not informed by the marker status. **b**, A non-targeted biomarker study in which the trial is designed and powered to address the biomarker hypothesis to ensure adequate biomarker

weak biological basis for differences that can be attributed to these superficial characteristics.

The recognition that clinical trials need to be redesigned to account for non-random variation comes more from knowledge of the disrupted cancer genome rather than of the germ line. The implications of having multiple potential treatments and diseases where once there was just one put enormous pressure on researchers to alter the design of clinical trials. Investigators often approach the challenge of having too many diseases and too few trial subjects as a result of genomic partitioning as a clinical-trial design problem. This creates unhealthy tension between design strategies because although clinical-trial design must be tailored to answer specific questions that arise from targeted therapies, many of these questions are actually standard and can be addressed by well-established methodologies. Consequently, the challenges of conducting clinical testing for most precision-medicine strategies revolve around their feasibility, efficiency and capacity to deal with multiple small-incidence subtypes of cancer and a rapidly evolving knowledge base.

In response, drug-development pathways have evolved to accommodate two important strategies: generating signals that indicate clearly the safety and efficacy of useful treatments, and terminating the development of ineffective treatments as early as possible. The four phases of clinical trials feed into these strategies. The early development phase (phase I) focuses on the safety aspects of a drug, including dosage, in a small group of patients. The middle-development phase (phase II) evaluates the safety and efficacy of a drug in a larger group of patients, and enables a 'go/no-go' decision to be

representation and distribution between arms. **c**, Biomarker-targeted randomized controlled trial (RCT) in which the presence of the selection marker guides patient allocation. **d**, RCT that compares biomarker-directed therapy with conventional therapy, which allows the overall concept of the biomarker approach to be tested as a whole. Adapted with permission from ref. 26. R, randomization.

made. The late development phase (phase III) constitutes comparative testing and provides a basis for seeking approval to market the drug. Phase IV trials are sometimes performed after market approval has been granted to examine the safety and efficacy of the drug in other patient populations, as well as any side effects and the implications of long-term use. These studies can also extend the applications or 'indications' of the drug. Through the sequential building of evidence, the use of a new therapeutic agent for a specific indication can be supported or refuted. In this model, a premium is placed on randomized, controlled designs.

Biomarkers — biological characteristics that can be measured in the context of diagnosis and clinical intervention — are often used to drive the selection of participants for trials, a strategy known as enrichment, which is well established for high-prevalence biomarkers. There are a number of methods for assessing the clinical utility of biomarkers (Fig. 1). For example, randomized controlled trial data can be analysed retrospectively (Fig. 1a). Biomarker discovery can also be integrated within the design of the trial to ensure that there is sufficient power to detect signals. Biomarker-positive patients can be equally distributed in each arm (known as biomarker stratification) to ensure statistical power (Fig. 1b), and the biomarker itself can be used to direct the study (Fig. 1c, d)^{26–28}. Advances in our understanding of the differences between the molecular pathologies of individual cancers creates challenges for conventional drug-development models, especially as the prevalence of molecular segments decreases²⁹. The chances of showing a significant effect in a traditional comparative trial of unselected participants diminish if the prevalence of a

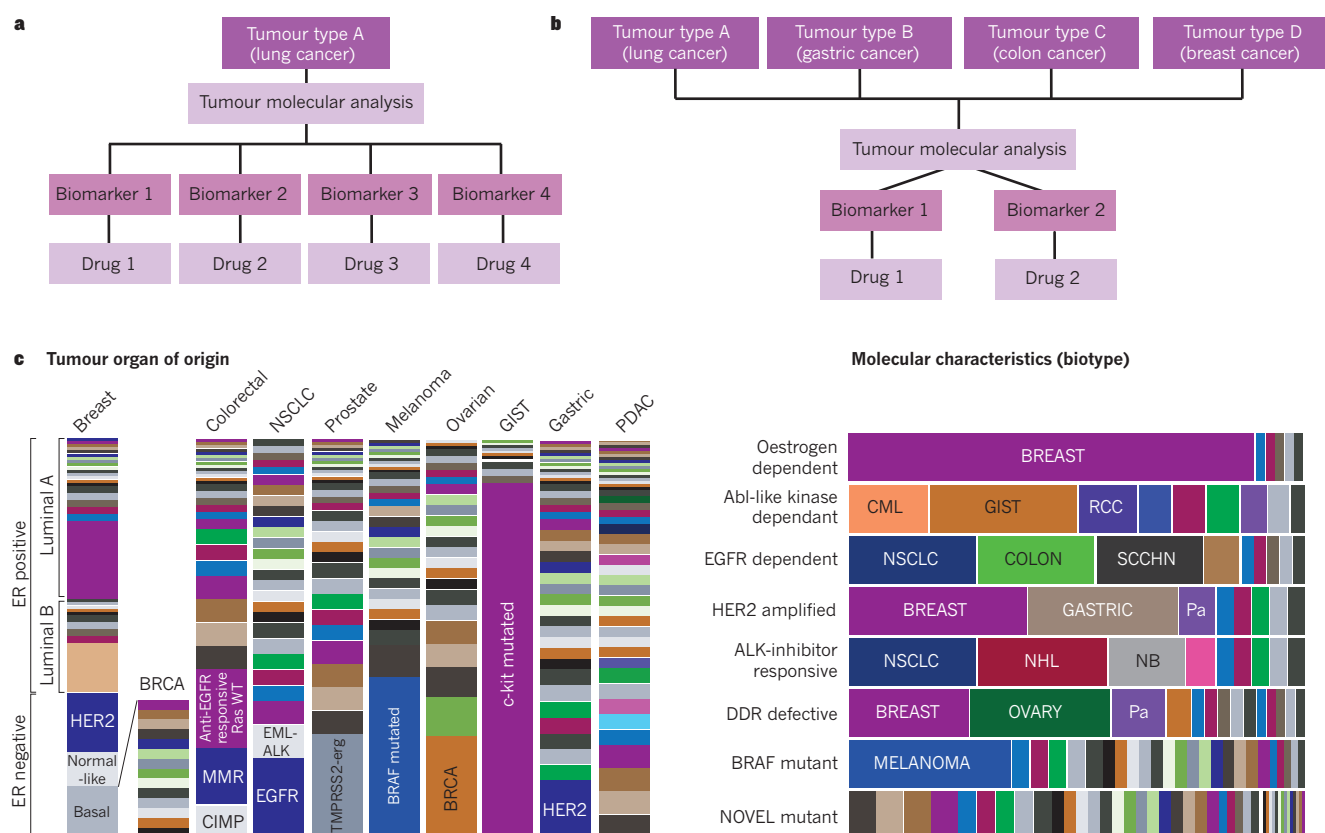


Figure 2 | Design principles that generate efficiencies in clinical trials of targeted therapies. **a**, In umbrella studies, patients with the same type of cancer are screened for a series of hypothesized predictive biomarkers. They are then allocated to appropriate therapies within the trial architecture. (The biomarker status for each tumour in the study is determined by tumour molecular analysis.) **b**, Basket studies recruit patients on the basis of their molecular characteristics irrespective of the organ in which their tumour originated. **c**, The relative incidence of molecular subtypes can help to guide decisions as to whether an umbrella or basket clinical-trial strategy is most appropriate. Molecular subtypes

can be classified by their organ of origin (left) or on the basis of their molecular characteristics or 'biotype' (right). Stratification is helpful when the incidence of a specific molecular class is low across different organs of origin and tends to be tested with a basket approach. Adapted with permission from ref. 75. CML, chronic myeloid leukaemia; DDR, DNA damage response; ER, oestrogen receptor; GIST, gastrointestinal stromal tumour; NB, neuroblastoma; NHL, non-Hodgkin lymphoma; NSCLC, non-small-cell lung cancer; Pa, pancreatic cancer; PDAC, pancreatic ductal adenocarcinoma (pancreatic cancer); RCC, renal cell carcinoma; SCCHN, squamous cell carcinoma of the head and neck; WT, wild type.

biomarker that identifies tumours most likely to respond to a targeted therapeutic agent is low. For example, if the biomarker is present in only 2% of the population — a typical prevalence for many, if not most, molecular segments³⁰ — a study of 50–100 patients yields only one or two patients. Unfortunately, no amount of clinical effect in such a small number of patients would be enough to advance the drug's therapeutic development (assuming that there is no clinical effect in the population who test negative for the biomarker).

Evaluating a targeted drug or treatment in the early phases of development will now more frequently require a trial with a selected patient population to minimize the inclusion of individuals who are unlikely to respond for mechanistic reasons. Inevitably, this yields smaller trials and fewer data on which to base decisions about trial-phase transitions. It also creates challenges when developing appropriate comparator populations in early studies. These approaches raise a number of interesting questions. For instance, how many patients must be evaluated to truly understand the safety and efficacy of a drug or treatment? Should later studies remain solely focused on the selected patient population and include just one arm? What are the drug effects in biomarker-negative patients? Owing to errors in diagnosis during routine clinical practice, such patient populations will exist even if they are not selected for investigation during the drug-development process. How can we build the body of evidence needed to support the approved use of a drug or therapeutic agent in a particular indication? As a consequence, challenges are introduced

throughout the entire drug-development pathway. These can be basic, such as the practicalities of finding enough patients who have low-incidence markers to investigate, and understanding the utility of the markers used for selection. They can also affect central aspects of the drug-development pathway, such as how to generate the data packages needed for regulatory submissions and market approval.

Patient-centric drug development

The challenges discussed in this Review have resulted in new clinical-trial designs (Fig. 2). An umbrella study (Fig. 2a) typically investigates a single tumour type selected according to the biomarkers relevant to one or more of the candidate drugs, and patients are directed towards different arms of the study — and hence towards different therapeutics — according to the molecular characteristics of their tumour. A basket study (Fig. 2b) also selects tumours according to their molecular characteristics and biomarkers, but is conducted irrespective of tumour type and often focuses on one (or a few) specific markers. The approach that is chosen will be based on various aspects, including the prevalence of a molecular subtype within a cancer type compared with its prevalence across different cancer types (Fig. 2c). Consideration will also be given to whether initiatives led by cooperative groups focussed on specific cancers exist, as well as the practicality of implementing these studies, such as the ability to acquire samples of tumour for analysis.

A solution to some of these challenges in targeted-drug

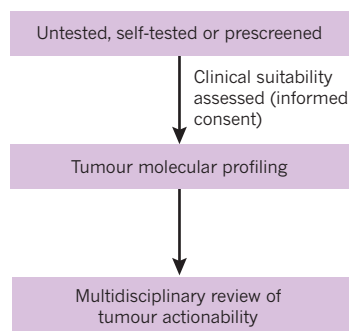
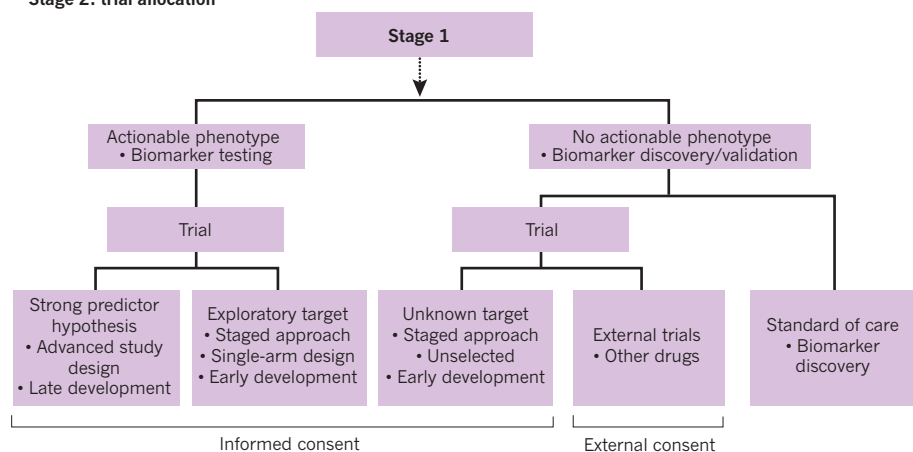
Stage 1: patient recruitment and molecular profiling**Stage 2: trial allocation**

Figure 3 | Master protocols for therapeutic development: a framework for the clinical testing of precision-oncology strategies — or ‘finding the trial for the patient’. A guiding principle within the framework is that all patients who are eligible for treatment should receive a choice of therapies. These therapies range from biomarker-directed or unselected new therapeutic strategies (either as part of the trial design or through external trials) to standard-of-care treatment in which patients will still be tracked to inform biomarker discovery opportunities for existing approved therapeutics. The framework can be enacted by a single

body or, more pragmatically, through a composite or network of organizations and activities with a co-ordinated management and governance structure. Stage 1 of the framework includes patient recruitment and molecular-testing. Participants are either screened before entering the trial or directed to molecular testing to be done within the trial structure itself or by external providers, if more appropriate. In stage 2, patients and clinicians are presented with a series of attractive clinical-trial options to choose from. This stage also incorporates an additional consent process.

development is the use of a master protocol, some of which have been established for efficiency in certain settings (Fig. 3 and Table 1). Rather than using serial, single diagnostic tests to select participants for different trials, a single, multiplex diagnostic assay is often used to assign participants to different candidate drugs (or arms of a trial) within the same trial, or a network of trials. This is sometimes referred to as a ‘tent’ protocol, in which multiple trials can be accessed through various mechanisms. Such studies offer more options for patients and can also make patient screening and recruitment more efficient.

Increasingly, adaptive design features are being incorporated. These differ from conventional designs by using accumulated results to modify the course or structure of a trial. The ability to make an early assessment of the clinical benefit or safety of a drug — and to modify the trial in response — is a nimble approach and offers a number of advantages. For instance, the trial can be stopped early or extended depending on the emerging results, or arms or doses can be dropped if no benefit is seen. This approach makes it easier to identify populations of patients who are responding to the drug being investigated, or to identify fruitful combinations of biomarkers and drugs or other therapeutics. It also allows the randomization proportions of the trial population or the rates at which data are accrued to be changed. Finally, it permits the inclusion of multiple stages of drug development within a single trial. Staged approaches such as these can markedly enable the drug-development process (Fig. 4). Examples of clinical trials that use these approaches include the Biomarker-integrated Approaches of Targeted Therapy for Lung Cancer Elimination (BATTLE)³¹ and the Investigation of Serial Studies to Predict Your Therapeutic Response with Imaging and Molecular Analysis (I-SPY) series^{32–36} of trials for lung and breast cancer.

Targeted therapeutic development is evolving rapidly, and there has been a notable expansion of precision-medicine programmes in recent years (Table 1). Combining a detailed understanding of the molecular pathology of tumours with modern drugs and associated diagnostic technologies for selecting patients has already translated into tangible improvements in survival rates for patients with certain cancer types^{10–17}, particularly those with non-small cell lung cancer (NSCLC)^{13,16}. In addition, significant durable responses to immune modulatory therapies have been discovered in about 15% of patients. These therapeutic agents target specific molecular mechanisms that

are currently the focus of intense investigation. Patient selection is also likely to play an important part in the development of these agents, with biomarker hypotheses being actively developed for the identification of trial participants³⁷. Data are emerging from early programmes such as SHIVA³⁸, which broadly evaluated targeted therapies without taking into account the histology of the tumour in end-stage patients for whom standard therapy had failed. Although no difference was identified³⁹, it is not possible to draw broad conclusions from this finding, exemplifying the challenges ahead.

The oncology landscape is accumulating a growing number of patient and tumour groups⁴⁰ that can be identified by (increasingly complex) diagnostic assays, which enables them to be coupled to molecularly targeted drugs. Up-to-date approvals can be found on the websites of the US Food and Drug Administration (FDA)^{41–43} and the European Medicines Agency (EMA)⁴⁴. Although most approved therapies have a linear relationship with a single biomarker, emerging data suggest that combinations of biomarkers might better inform therapeutic responsiveness, and will continue to challenge biomarker development. Similarly, multiple biomarkers could indicate sensitivity to a single therapeutic agent, and conversely a single biomarker might define patients that would benefit from several therapeutic options. Such overlaps are inevitable and it is important to define appropriate measures on how to respond to them during the drug-development process. The emerging complexity poses substantial challenges for current regulatory processes. For example, how should researchers assess therapies that do not take the cancer’s organ of origin into account, particularly when its prevalence is low in a particular organ? How should therapies be assessed at different stages of the disease, especially in cases where the patient has undergone several prior treatments? A solution might be to apply a broader approach, such as defining the level of reimbursement for a particular disease stage and line of treatment, with decisions on choice of therapy made between clinicians and their patients.

The challenges of early drug development

Clinical testing in the early stage of drug development poorly predicts efficacy in later stages of development^{25,45}. Bias in small early trials can raise expectations, only to cause disappointment when they are expanded to include larger, less-selected and unbiased populations. Current tools that provide an improved understanding of the

Table 1 | Precision-medicine studies

Precision-medicine clinical trials							
Study	Tumour	Phase/design	Location	Arms	Patients†	Clinical trial ID	References
Bisgrove	All	Phase II, non-randomized	United States	N/A	84	NCT00530192	19
IMPACT	All	Phase I	United States	N/A	1,144	NCT00851032	20
MOSCATO 01	All	Phase I	France	N/A	420	NCT01566019	21
Lung-MAP	Squamous lung	Phase II/III, randomized	United States	5	10,000	NCT02154490	49
BATTLE	NSCLC	Umbrella, route to four phase II randomized	United States	4	300	NCT00409968 (umbrella) NCT00411671 NCT00411632 NCT00410059 NCT00410189	31, 66, 67
BATTLE-2	NSCLC	Phase II randomized	United States	4	450	NCT01248247	N/A
BATTLE-FL	NSCLC	Phase II randomized	United States	4	225	NCT01263782	N/A
I-SPY 2	Breast cancer	Phase II randomized	United States	8	800	NCT01042379	68, 69
NCI-IMPACT	All	Phase II stratified, non-randomized	United States	6	700	NCT01827384	70
NCI-MATCH	Solid	Phase II stratified, non-randomized	United States	20	3,000	Umbrella, route to phase II‡	48
V-BASKET	All	Phase II stratified, non-randomized	Global	2	160	NCT01524978	71
CREATE	Selected	Phase II stratified, non-randomized	European Union	6	582	NCT01524926	N/A
WINTHER	All	Stratified, non-randomized	European Union	2	200	NCT01856296	72
SHIVA	All	Phase II stratified, controlled	France	10	1,000	NCT01771458	38
MOST	All	Phase II stratified, randomized	France	5	560	NCT02029001	N/A
SAFIR 02 Lung	NSCLC	Phase II stratified, randomized	France	8	650	NCT02117167	73
SAFIR 02 Breast	Breast cancer	Phase II stratified, randomized	France	18	460	NCT02299999	N/A
Lung MATRIX	NSCLC	Phase II stratified, non-randomized	United Kingdom	21§	2,000	EudraCT 2014-000814-73	65
FOCUS 4	Colorectal cancer	Phase II/III randomized	United Kingdom	4	643	EudraCT 2012-005111-12	74
IMPACT	Pancreatic cancer	Phase II stratified, randomized	Australia	4	90	ACTRN 12612000777897	47
Screening programmes that feed into precision-medicine trials							
Study	Tumour	Phase/design	Location	Diagnostics	Patients†	Clinical trial ID	References
I-SPY	Breast cancer	Phase II, diagnostic study	United States	Genomic, imaging	221	NCT00043017	32–35
NCI-MATCH	Solid	Screening, route to phase II	United States	NGS¶	3,000	N/A	48
VIKTORY	Gastric cancer	Screening, route to phase II	Asia	NGS, other#	600	NCT02299648	N/A
LC-SCRUM	NSCLC	Screening, route to phase II/III	Asia	As needed**	Open††	N/A	53
AURORA	Breast cancer	Screening, route to phase I/II/III	European Union	NGS, other‡‡	1,300	NCT02102165	52
SPECTAColor	Colorectal cancer	Screening, route to phase I/II/III	European Union	NGS	2,600	NCT01723969	50
SPECTALung	Lung	Screening, route to phase I/II/III	European Union	NGS	500§§	NCT02214134	51
MOSCATO	All	Screening, route to phase I/II	France	CGH array, sequencing	1,050	NCT01566019	21
SAFIR 01	Breast cancer	Screening, route to phase I/II	France	CGH, sequencing, gene expression array	423	NCT01414933	73
CRUK SMP1	Selected	Screening, feasibility	United Kingdom	Bespoke panel	9,000	N/A	36

BATTLE-FL, Front-Line Biomarker-Integrated Treatment Study in Non Small Cell Lung Cancer; CGH, comparative genomic hybridization; FISH, fluorescence *in situ* hybridization; IHC, immunohistochemistry; IMPACT, Individualised Molecular Pancreatic Cancer Therapy; IMPACT, Initiative for Molecular Profiling in Advanced Cancer Therapy; MOSCATO, Molecular Screening for Cancer Treatment Optimization; MOST, Adapting Treatment to the Tumor Molecular Alterations for Patients with Advanced Solid Tumors: My Own Specific Treatment; N/A, not applicable; NCI-IMPACT, National Cancer Institute-Molecular Profiling-Based Assignment of Cancer Therapy for Patients with Advanced Solid Tumors; NGS, next-generation sequencing; VIKTORY, Targeted Agent Evaluation in Gastric Cancer Basket Korea Study.

†Estimated number of patients to be recruited, or the final number recruited where the study has been completed. ‡The NCI-MATCH programme is a screening programme used to direct patients to single-arm, phase II, signal-seeking studies. §The number of arms will vary because the study progresses as each arm has been designed around a biomarker (for patient selection) and (candidate) drug pair. ||Once fully operational, the study will screen 2,000 patients per year. ¶FISH and IHC assays will be used as required. #‘Other’ refers to a selection of bespoke and exploratory diagnostics. **Bespoke diagnostics are deployed as needed to select patients for the individual clinical studies that feed from the screening programme. ††‘Open’ describes an open and rolling patient-recruitment programme. ‡‡‘Other’ refers to RNA sequencing. §§500 patients in year 1 then 500–1000 patients, thereafter.

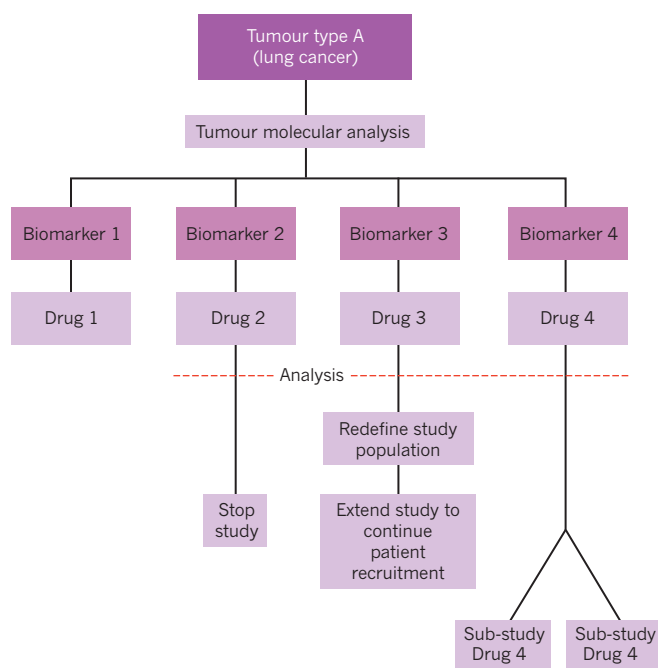


Figure 4 | Adaptive study designs. A within-study analysis or the continual assessment of data can be used to change the course of a clinical trial. First, the biomarker status for each tumour in the study is determined by tumour molecular analysis. After each tumour is allocated to a suitable sub-study, further analysis is conducted. Consequently, the sub-study 2 trial arm can be stopped owing to a lack of evidence to support the clinical benefit of drug 2, and the sub-study 3 trial arm can be extended to include more patients. Meanwhile, the patient population of sub-study 4 can be redefined into two sub-studies, according to the results of responder/non-responder analysis.

molecular pathology of tumours can be used to inform smaller trials as well as to define sources of bias at the molecular level to inform early and ongoing therapeutic development. An emerging approach

is the testing of small numbers of patients underpinned by a deep understanding of both the molecular composition of tumours and the mechanism of action of the therapeutic agent. Knowledge acquired through clinical testing can then inform ongoing preclinical strategies, which in turn refine the clinical-testing approach — a process known as forward-and-backward translation (Fig. 5).

Inherent to this approach is a desire to define more effective therapies and to set the bar higher for furthering the progression of a therapeutic agent down the drug-development pathway. A shift is needed away from the current high-investment drug-development approach that is dominated by late-phase trials that predominantly fail at great expense, towards an approach in which failures are early and cheap. This will allow a greater number of potential therapies to be assessed while constraining costs. Researchers might even be able to test bolder biological hypotheses, particularly in cancers for which current therapeutic options are poor. With these tools in hand, and developing rapidly, the challenge now becomes to determine how we can implement these strategies in the real world.

Master-protocol clinical trials that use umbrella and basket designs to enable trial stages to be run in parallel are efficient. However, the subdivision of tumour and therapeutic pairs that they create highlights a need for more innovative solutions and approaches, particularly in early drug development^{27,46}. For example, there might not be enough patients to test the targeted therapeutic using conventional designs. Figure 6 shows a suggested strategy for the development of therapeutic agents to treat cancer with an overall incidence of 10 patients per 100,000 individuals per year. Supportive evidence for a particular strategy can be classified according to an ‘actionability index’. The development of each therapeutic agent will progress within this framework or graduate to pivotal studies when there is sufficient evidence.

Accelerating stratified therapeutic development

The development of precision therapeutics focuses on leveraging the science, however, many important challenges pivot on operational components⁴⁷. These components require the integration of multiple complex processes such as participant screening and recruitment

BOX 1

Delivering multidrug–portfolio studies

A number of diagnostic, protocol and operational requirements must be considered when designing clinical trials that use multidrug portfolios.

● **Participant screening and recruitment** There should be a viable means by which to identify low-incidence patient subpopulations and to direct individuals to an appropriate clinical trial. Patient-centric approaches give individuals access to many options through a single screening process. Such screening programmes are usually region-wide and collaborative. They can be linked to umbrella and basket studies and also to global studies that accept participants from diverse screening routes. Drug portfolios are made available to these trials through collaborations, and safeguards are implemented for proprietary information when multiple partners are involved. The multiplexed diagnostic platforms and systems should be harmonized or cross-validated to allow patients to be recruited irrespective of the technology used by partners. Regulators should be open to changes with respect to how these clinical trials are run. The screening programmes are underpinned by networks, collaborations and reliable partners.

● **Molecular testing** The testing platform and screening or selection algorithm should enable broad yet robust tumour and patient profiling.

They should provide viable drug-development routes for larger or global studies, regulatory interactions and markets. Samples must be used efficiently and data generation should be robust. Overall, molecular tests should be cost-effective, transferable and widely deployable. Testing should be performed to agreed standards.

● **Protocols** Trials should start with a flexible protocol that can incorporate both emerging changes in the science and an understanding of patient and tumour biomarkers. Alternatively, they could use a confirmatory development protocol that permits regulatory interactions that accept different types of data. Such protocols can be deployed on their own or in alignment with other protocols. They can be modular, rolling or open ended, and must be reviewed efficiently according to a centralized regulatory and ethics process.

● **Availability and delivery of therapies** Operational machinery must be chosen that allows clinical studies to be conducted in diverse groups of patients and over a broad geographical area. Regulatory and ethics processes and patient screening and recruitment should be aligned and efficient. Therapies can be distributed using hub-and-spoke models and cost-effective and efficient delivery of multiple candidate drugs to multiple sites can be facilitated through a centralized pharmacy. The work should be highly collaborative, spread across many groups and involve reliable partners.

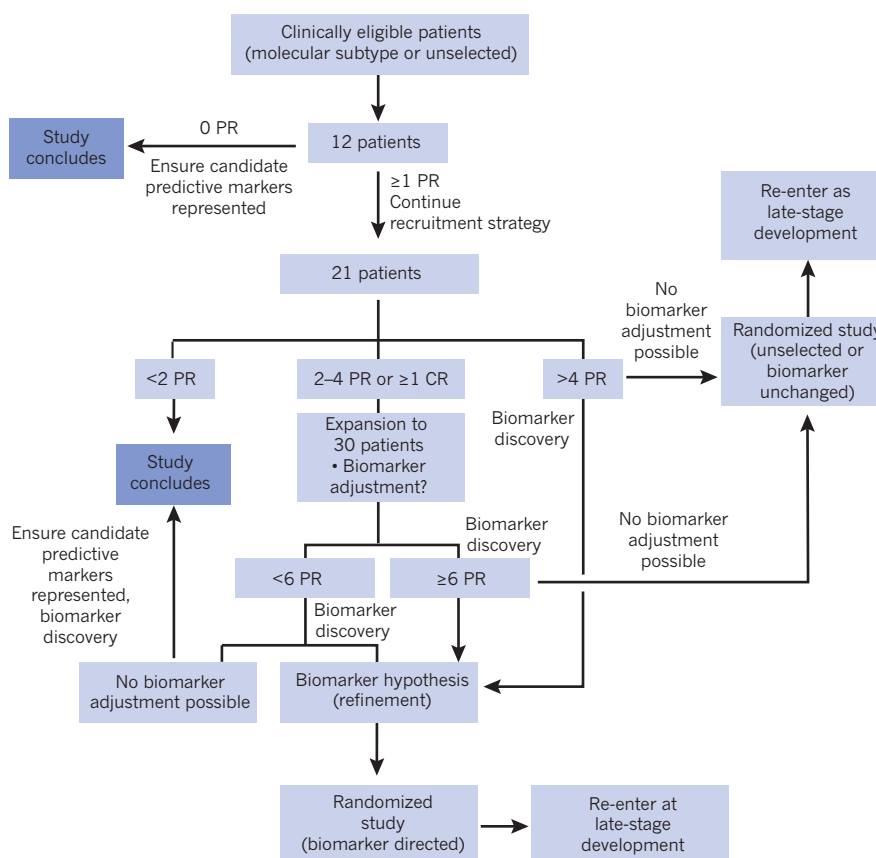


Figure 5 | Early stratified therapeutic development. An important element of early therapeutic development is the use of small trials that are underpinned by a deep understanding of tumour molecular pathology, which guides ongoing trial development. A stepwise development approach is applied, and

interim analyses, trial-population expansions and molecular assessments are implemented at specific points. CR, complete response, and PR, partial response, based on Response Evaluation Criteria in Solid Tumors (RECIST) 1.0 criteria⁷⁶.

and the molecular testing of tumours. Rather than pursuing the conventional goal of finding the patient for the trial, the overall goal is to 'find the trial for the patient'. Protocols must also be flexible and therapies must be available and deliverable (Box 1).

Participant screening and recruitment

The realities of the conventional screening approach in clinical drug development are sobering. For example, consider a candidate-drug trial in a subpopulation of patients that were selected by a biomarker with a 2% incidence, which has a typical screening failure rate of 15% and a patient dropout rate of 15%. The trial would need to screen 78 patients to find one patient for recruitment, which effectively means that 77 patients are discarded. The cost of such an approach is equally sobering. Screening using routine single-variable diagnostic approaches, such as immunohistochemistry or a single-gene DNA test, would have a cost of about US\$1,125 per assay, which includes performing and processing the assay, as well as logistics and reporting. It would therefore cost \$88,235 to screen enough individuals to recruit one participant. To conduct a 20-patient phase I expansion study in this selected patient subpopulation, the trial would need to screen 1,560 patients, at a cost of \$1.8 million.

In addition, the patient's experience during the conventional screening approach is often extremely poor and can involve many cycles of disappointment. After first being considered for a trial, the patient might then become ineligible to participate if they do not have the correct biomarker. They must then undergo repeat biopsies during the search for the next biomarker, and ultimately might receive only limited drug options. The physician's experience is similarly poor: his or her options are limited to screening for different biomarkers, and associated trials, so long as tumour material is available.

From the operational viewpoint of a clinical trial, this is unsustainable for practical reasons, such as the lack of available tissue and the unwillingness of patients and clinicians to participate.

The need to find sufficient numbers of patients with a specific biomarker has generated many cooperative study groups (Table 1). Consortia provide multiplexed molecular testing assays — in which many biomarkers are measured concurrently — as part of the drug-development process, as well as programmes that offer 'self-tested' patients access to appropriate therapy either as part of clinical trials or through 'off-label' treatment. In the United States, examples include national-level, cross-sector collaborative (including government-based) initiatives such as the National Cancer Institute-Molecular Analysis for Therapy Choice (NCI-MATCH)⁴⁸ (solid tumours) and Lung Cancer Master Protocol (Lung-MAP, NCT number NCT02154490)⁴⁹ (squamous lung cancer) programmes. Other examples include the Screening Patients for Efficient Clinical Trial Access (SPECTA) programmes (SPECTAColor⁵⁰ in colorectal cancer (NCT01723969) and SPECTALung⁵¹ in lung cancer (NCT02214134)) and the AURORA initiative in Europe⁵² (breast cancer (NCT02102165)), and the Lung Cancer Genomic Screening Project for Individualized Medicine in Japan (LC-SCRUM-Japan)⁵³. Cancer-specific advocacy and charity organizations also lead cooperative study groups, such as the 'Know Your Tumor' programme established by the Pancreatic Cancer Action Network in the United States. Although these models are advancing precision oncology, they are costly because they require intermediaries to navigate the patient through the health-care system. They are also difficult to scale up without fundamental changes in health-service delivery. Meanwhile, patients and clinicians are also driving forwards new approaches. These approaches include clinical trials and other therapeutic options as part of a molecular assay report,

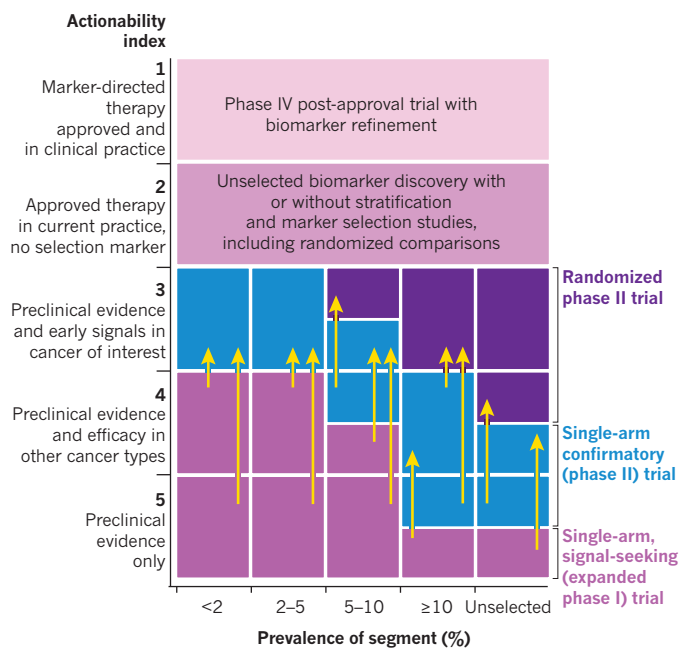


Figure 6 | Clinical-testing strategies. Lower-prevalence segments present a considerable challenge when testing stratified therapeutic strategies. It is also a challenge to determine the level of evidence that is required to embark on later-phase studies. The potential approach shown in this matrix is a function of the existing level of evidence, the prevalence of the segment, which indicates the feasibility of the testing strategy, and current regulatory requirements. Trials can progress as the level of evidence increases, and this progression can be built into the planned stepwise development process.

such as FoundationOne from Foundation Medicine, and connections to further information, consumer-focused advice, communities and patient-led consortia. The broader net that such approaches cast helps to identify smaller and smaller subtypes and opportunities for individual patients. Strategies that provide genomic health advice⁵⁴ and navigation, such as Perthera, are also gaining traction. Others have begun to use electronic media to enable patients and clinicians to ‘shop around’ for the best option. These strategies can markedly improve efficiency and the patient experience. However, despite these efforts, trials using a selection biomarker still constitute only a minority of current studies⁵⁵.

Recruiting eligible patients onto a clinical trial represents a major challenge. If the prevalence of eligible patients is low, it is often necessary to open a large number of screening centres — a considerable cost, especially since not all will be able to recruit patients. As screening programmes expand in size, the cost of funding the search for patients shifts from drug developers to health-care systems or research platforms. A possible solution is to open clinical trials at a location that is accessible to the patient only after they have been identified — known as ‘just-in-time’ accessibility. The cost of rapidly deploying teams to establish a trial location after a patient has been found is likely to be lower than the cost of screening a large number of patients.

Molecular testing

Although multiplex testing of the coding regions of candidate genes offers some options, the complexity of cancer will inevitably require more in-depth analyses⁵⁶. The challenges of delivering molecular assays using advanced technologies are discussed elsewhere⁵⁶, however, current tests exploit the relatively direct relationships that exist between a specific mutation and the efficacy of a drug. The appraisal and delivery of more complex assays that might better identify responsive subtypes^{57,58} is proving to be difficult despite advances in clinical-grade diagnostics⁵⁹. This is mainly due to the rigidity and

inertia of established processes for biospecimen handling. Simple solutions such as liquid biopsies^{60,61} are promising, but could lack broad applicability, particularly when complex molecular changes must be analysed. Technology considerations aside, it is more important to understand the relevance of any detected changes or mutations, and the body of evidence that is required to substantiate their use for patient selection. Modern multiplex systems such as next-generation sequencing technologies reveal the molecular changes within a single tumour at an unprecedented level of detail. Many of these changes will not have been widely reported: some are likely to be specific to that tumour (or tumour region) and there will be little previous clinical experience or knowledge for most. In light of this, how should therapeutic selection be informed? Although specific mutations in a particular gene can confer sensitivity to a particular therapeutic agent, what should we do if we discover previously unreported mutations in that same gene? And what should we do if the potential functional consequences have not been investigated yet? Can these mutations reasonably be expected to confer similar therapeutic sensitivity? This challenge is being addressed through trial design and the diagnostic algorithms that are used to assign patients to treatments. We must be careful to avoid reporting a study as negative purely because it has not shown any clinical benefit in a subpopulation that has been defined by mutations of unknown consequence. Not all mutations in a gene will be predictive of clinical benefit. Practical solutions to accommodate such uncertainty often combine adaptations within umbrella- or basket-shaped trial arms that can examine combinations of biomarkers and therapeutics in isolation. Different weightings can then be attributed to mutations of known and unknown clinical or functional consequence — a process called mutation tiering in which groups are designated as either ‘tight’ markers that have a high level of supportive evidence or ‘loose’ markers that are more exploratory in nature.

Protocol flexibility

The administrative and logistical challenges of clinical trials are substantial. They impede the ability to respond nimbly to trial findings, particularly if unexpected, or to data emerging from outside the trial. Establishing frameworks and platforms for stratified therapeutics development will facilitate the deployment of ‘within-protocol’ responses to specific scenarios, which will improve flexibility of trials (Box 1).

Availability and delivery of therapeutics

Conducting molecular analysis without the prospect of a resulting action is of little value. There are comparatively few opportunities in routine health care in which multiplexed testing can be applied to influence clinical decision-making, and access to appropriate therapeutics remains problematic⁶². Negotiating individual clinical trials on an *ad hoc* basis is impractical because of slow legal and administrative processes — a closer relationship must be cultivated between the pharmaceutical industry and other stakeholders to ease this roadblock. The involvement of multiple pharmaceutical partners will ensure that a broader range of candidate drugs and appropriate comparator therapies are available. Wider collaboration between tumour-specific consortia, diagnostic and regulatory groups, as well as major charities and other interested parties, will also be pivotal. A drug-portfolio approach — negotiated as a broad partnership or through a consortium strategy — is a necessity, as is the ability to deliver therapeutic agents through systems such as a centralized pharmacy. The ability to offer patients and clinicians a broad selection of attractive treatment options will enhance participation in clinical trials. At present, only 2–5% of potentially eligible participants^{63,64} enrol in such trials. Initiatives such as NCI-MATCH⁴⁷, Lung-MAP⁴⁹, and the Cancer Research UK Stratified Medicine programmes and the National Lung Matrix Trial (Lung MATRIX)⁶⁵ (European Clinical Trials Database (EudraCT) number 2014-000814-73) have set a

precedent for prioritizing participation rates. However, the real value to the patient and health-care system will be when these strategies become commonplace and encompass a greater proportion of drug-development portfolios. This will ensure the broad availability of therapeutics currently in development.

Most advances have been achieved by altering drug-development strategies to fit into established health-care systems. Consequently, progress has been slow. If health-care systems are out of pace with the drug-development process, they could be impeding the development of therapeutic agents. Health-care systems that can implement precision medicine will greatly facilitate therapeutic development. To accelerate progress, health-care systems must be aligned to ensure that they are able to test and deliver precision medicine without the need for costly overlying clinical-trial infrastructure.

Future directions

In recent years, our understanding of the precision-therapeutic development pathway has evolved rapidly. In some areas, targeted-drug development has progressed from concept to reality. The frameworks, platforms and processes involved are now capable of supporting modern oncology drug development. Innovative clinical-trial designs — also a central component of development — are highlighting the need to better appraise tumour biology, drug efficacy and the potential benefits for patients. Emerging drug-development paradigms are driving new ways of working collaboratively to accelerate progress. By generating truly patient-centric clinical trials, we have taken important early steps into the evolving era of precision medicine. In some cases, these steps are already enabling us to 'select' the trial for the patient. However, major hurdles remain, and we must establish broad frameworks and systems that integrate closely with health-care delivery to accelerate progress and realize the true promise of precision medicine. ■

Received 20 May; accepted 14 August 2015.

1. Chin, L. & Gray, J. W. Translating insights from the cancer genome into clinical practice. *Nature* **452**, 553–563 (2008).
A review that outlines the opportunities, challenges and approaches associated with the advancement of genomics-based medicine.
2. Stratton, M. R. Exploring the genomes of cancer cells: progress and promise. *Science* **331**, 1553–1558 (2011).
3. Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. *Nature* **458**, 719–724 (2009).
A review of recent progress in cancer genomics and the potential of its application to medicine.
4. Hudson, T. J. *et al.* International network of cancer genome projects. *Nature* **464**, 993–998 (2010); erratum **465**, 966 (2010).
5. The Cancer Genome Atlas Research Network. Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature* **455**, 1061–1068 (2008); erratum **494**, 506 (2013).
6. Chin, L., Andersen, J. N. & Futreal, P. A. Cancer genomics: from discovery science to personalized medicine. *Nature Med.* **17**, 297–303 (2011).
A review that addresses the accumulating knowledge acquired through large-scale genomic sequencing efforts and discusses strategies for translating these discoveries into patient care.
7. Zhang, B. *et al.* Proteogenomic characterization of human colon and rectal cancer. *Nature* **513**, 382–387 (2014).
8. Verweij, J. *et al.* Progression-free survival in gastrointestinal stromal tumours with high-dose imatinib: randomised trial. *Lancet* **364**, 1127–1134 (2004).
9. Gerber, D. E. & Minna, J. D. ALK inhibition for non-small cell lung cancer: from discovery to therapy in record time. *Cancer Cell* **18**, 548–551 (2010).
10. Sosman, J. A. *et al.* Survival in BRAF V600-mutant advanced melanoma treated with vemurafenib. *N. Engl. J. Med.* **366**, 707–714 (2012).
11. Slamon, D. *et al.* Adjuvant trastuzumab in HER2-positive breast cancer. *N. Engl. J. Med.* **365**, 1273–1283 (2011).
12. Shaw, A. T. *et al.* Crizotinib versus chemotherapy in advanced ALK-positive lung cancer. *N. Engl. J. Med.* **368**, 2385–2394 (2013).
13. Maemondo, M. *et al.* Gefitinib or chemotherapy for non-small-cell lung cancer with mutated EGFR. *N. Engl. J. Med.* **362**, 2380–2388 (2010).
14. Ledermann, J. *et al.* Olaparib maintenance therapy in patients with platinum-sensitive relapsed serous ovarian cancer: a preplanned retrospective analysis of outcomes by BRCA status in a randomised phase 2 trial. *Lancet Oncol.* **15**, 852–861 (2014).
15. Kris, M. G. *et al.* Using multiplexed assays of oncogenic drivers in lung cancers to select targeted drugs. *J. Am. Med. Assoc.* **311**, 1998–2006 (2014).
16. Jänne, P. A. *et al.* AZD9291 in EGFR inhibitor-resistant non-small-cell lung cancer. *N. Engl. J. Med.* **372**, 1689–1699 (2015).
17. Demetri, G. D. *et al.* Efficacy and safety of imatinib mesylate in advanced gastrointestinal stromal tumors. *N. Engl. J. Med.* **347**, 472–480 (2002).
18. Green, E. D., Guyer, M. S. & National Human Genome Research Institute. Charting a course for genomic medicine from base pairs to bedside. *Nature* **470**, 204–213 (2011).
A perspective article that describes the past, present and future trajectories of genomic medicine.
19. Von Hoff, D. D. *et al.* Pilot study using molecular profiling of patients' tumors to find potential targets and select treatments for their refractory cancers. *J. Clin. Oncol.* **28**, 4877–4883 (2010).
This paper and refs 20 and 21 are some of the first descriptions of the use of molecular targeted therapies to improve patient outcomes.
20. Tsimberidou, A.-M. *et al.* Personalized medicine in a phase I clinical trials program: the MD Anderson Cancer Center initiative. *Clin. Cancer Res.* **18**, 6373–6383 (2012).
21. Hollebecque, A. *et al.* Molecular screening for cancer treatment optimization (MOSCATO 01): a prospective molecular triage trial — interim results. *J. Clin. Oncol.* **31**, 2512 (2013).
22. Dienstmann, R. *et al.* Molecular profiling of patients with colorectal cancer and matched targeted therapy in phase I clinical trials. *Mol. Cancer Ther.* **11**, 2062–2071 (2012).
23. Association of the British Pharmaceutical Industry. *The Stratification of Disease for Personalised Medicines* http://www.abpi.org.uk/our-work/library/medical-disease/Documents/strat_med.pdf (2014).
24. The Academy of Medical Sciences. *Realising the Potential of Stratified Medicine* <https://www.acmedsci.ac.uk/viewFile/51e915f9f09fb.pdf> (2013).
25. Cook, D. *et al.* Lessons learned from the fate of AstraZeneca's drug pipeline: a five-dimensional framework. *Nature Rev. Drug Discov.* **13**, 419–431 (2014).
26. Lee, C. K., Lord, S. J., Coates, A. S. & Simes, R. J. Molecular biomarkers to individualise treatment: assessing the evidence. *Med. J. Aust.* **190**, 631–636 (2009).
27. Sargent, D. J., Conley, B. A., Allegra, C. & Collette, L. Clinical trial designs for predictive marker validation in cancer treatment trials. *J. Clin. Oncol.* **23**, 2020–2027 (2005).
A description of the fundamental basis of clinical-trial designs that are used to assess biomarkers.
28. Mandrekar, S. J. & Sargent, D. J. Clinical trial designs for predictive biomarker validation: theoretical considerations and practical challenges. *J. Clin. Oncol.* **27**, 4027–4034 (2009).
29. Printz, C. Failure rate: why many cancer drugs don't receive FDA approval, and what can be done about it. *Cancer* **121**, 1529–1530 (2015).
30. Sleijfer, S., Bogaerts, J. & Siu, L. L. Designing transformative clinical trials in the cancer genome era. *J. Clin. Oncol.* **31**, 1834–1841 (2013).
31. Kim, E. S. *et al.* The BATTLE trial: personalizing therapy for lung cancer. *Cancer Discov.* **1**, 44–53 (2011).
32. Esserman, L. J. *et al.* Chemotherapy response and recurrence-free survival in neoadjuvant breast cancer depends on biomarker profiles: results from the I-SPY 1 TRIAL (CALGB 150007/150012; ACRIN 6657). *Breast Cancer Res. Treat.* **132**, 1049–1062 (2012).
33. Esserman, L. J. *et al.* Pathologic complete response predicts recurrence-free survival more effectively by cancer subset: results from the I-SPY 1 TRIAL–CALGB 150007/150012, ACRIN 6657. *J. Clin. Oncol.* **30**, 3242–3249 (2012).
34. Hylton, N. M. *et al.* Locally advanced breast cancer: MR imaging for prediction of response to neoadjuvant chemotherapy—results from ACRIN 6657/I-SPY TRIAL. *Radiology* **263**, 663–672 (2012).
35. Lin, C. *et al.* Locally advanced breast cancers are more likely to present as Interval Cancers: results from the I-SPY 1 TRIAL (CALGB 150007/150012, ACRIN 6657, InterSPORE Trial). *Breast Cancer Res. Treat.* **132**, 871–879 (2012).
36. Lindsay, C. R., Shaw, E., Walker, I. & Johnson, P. W. Lessons for molecular diagnostics in oncology from the Cancer Research UK Stratified Medicine Programme. *Expert Rev. Mol. Diagn.* **15**, 287–289 (2015).
37. Le, D. T. *et al.* PD-1 Blockade in tumors with mismatch-repair deficiency. *N. Engl. J. Med.* **372**, 2509–2520 (2015).
38. Le Tourneau, C. *et al.* Designs and challenges for personalized medicine studies in oncology: focus on the SHIVA trial. *Target. Oncol.* **7**, 253–265 (2012).
39. Le Tourneau, C. *et al.* Randomized phase II trial comparing molecularly targeted therapy based on tumor molecular profiling versus conventional therapy in patients with refractory cancer: results of the SHIVA trial. *J. Clin. Oncol.* **33**, 11113 (2015).
40. Watson, I. R., Takahashi, K., Futreal, P. A. & Chin, L. Emerging patterns of somatic mutations in cancer. *Nature Rev. Genet.* **14**, 703–718 (2013).
41. US Food and Drug Administration. Nucleic Acid Based Tests. *US Food and Drug Administration* <http://www.fda.gov/MedicalDevices/ProductsandMedicalProcedures/InVitroDiagnostics/ucm330711.htm> (2015).
42. US Food and Drug Administration. List of Cleared or Approved Companion Diagnostic Devices (In vitro and Imaging Tools). *US Food and Drug Administration* <http://www.fda.gov/MedicalDevices/ProductsandMedicalProcedures/InVitroDiagnostics/ucm301431.htm> (2015).
43. US Food and Drug Administration. Drug Approvals And Databases. *US Food*

- and Drug Administration <http://www.fda.gov/Drugs/InformationOnDrugs/> (2015).
44. European Medicines Agency. European public assessment reports. *European Medicines Agency* http://www.ema.europa.eu/ema/index.jsp?curl=pages/medicines/landing/epar_search.jsp&mid=WC0b01ac058001d125 (2015).
 45. Hay, M., Thomas, D. W., Craighead, J. L., Economides, C. & Rosenthal, J. Clinical development success rates for investigational drugs. *Nature Biotechnol.* **32**, 40–51 (2014).
 46. Yap, T. A., Sandhu, S. K., Workman, P. & de Bono, J. S. Envisioning the future of early anticancer drug development. *Nature Rev. Cancer* **10**, 514–523 (2010).
 47. Chantrill, L. A. *et al.* Precision medicine for advanced pancreas cancer: the Individualized Molecular Pancreatic Cancer Therapy (IMPACT) trial. *Clin. Cancer Res.* **21**, 2029–2037 (2015).
 48. National Cancer Institute *Molecular Analysis for Therapy Choice* http://deainfo.nci.nih.gov/advisory/ncab/164_1213/Conley.pdf.
 49. Lung Cancer Master Protocol (Lung-MAP) Clinical Trials. About Lung-MAP. *Lung-MAP* <http://www.lung-map.org/about-lung-map> (2015).
 50. EORTC. About SPECTAColor. *SPECTAColor EORTC Colorectal Cancer Screening Platform* <http://spectacolor.eortc.org/about> (2015).
 51. EORTC. EORTC, through SPECTALung, participates in EU consortium validating blood-based cancer biomarkers. *EORTC The future of cancer therapy* <http://www.eortc.org/news/eortc-through-spectalung-participates-in-european-consortium-validating-blood-based-cancer-biomarkers/> (2015).
 52. Zardavas, D. *et al.* The AURORA initiative for metastatic breast cancer. *Br. J. Cancer* **111**, 1881–1887 (2014).
 53. Matsumoto, S. *et al.* Nationwide genomic screening network for the development of novel targeted therapies in advanced non-small cell lung cancer (LC-SCRUM-Japan). *J. Clin. Oncol.* **33** (suppl. 15), 8093 (2015).
 54. Kalf, R. R. *et al.* Variations in predicted risks in personal genome testing for common complex diseases. *Genet. Med.* **16**, 85–91 (2014).
 55. Roper, N., Stensland, K. D., Hendricks, R. & Galsky, M. D. The landscape of precision cancer medicine clinical trials in the United States. *Cancer Treat. Rev.* **41**, 385–390 (2015).
 56. Simon, R. & Roychowdhury, S. Implementing personalized cancer genomics in clinical trials. *Nature Rev. Drug Discov.* **12**, 358–369 (2013).
 57. Waddell, N. *et al.* Whole genomes redefine the mutational landscape of pancreatic cancer. *Nature* **518**, 495–501 (2015).
- A report that demonstrates how different genomic readouts could be important biomarkers for therapeutic responsiveness.**
58. Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J. & Stratton, M. R. Deciphering signatures of mutational processes operative in human cancer. *Cell Rep.* **3**, 246–259 (2013).
 59. Frampton, G. M. *et al.* Development and validation of a clinical cancer genomic profiling test based on massively parallel DNA sequencing. *Nature Biotechnol.* **31**, 1023–1031 (2013).
 60. Dawson, S. J. *et al.* Analysis of circulating tumor DNA to monitor metastatic breast cancer. *N. Engl. J. Med.* **368**, 1199–1209 (2013).
 61. Douillard, J. Y. *et al.* Gefitinib treatment in *EGFR* mutated caucasian NSCLC: circulating-free tumor DNA as a surrogate for determination of *EGFR* status. *J. Thorac. Oncol.* **9**, 1345–1353 (2014).
 62. Lewin, J. & Siu, L. L. Cancer genomics: the challenge of drug accessibility. *Curr. Opin. Oncol.* **27**, 250–257 (2015).
 63. Lara, P. N. Jr *et al.* Prospective evaluation of cancer clinical trial accrual patterns: identifying potential barriers to enrollment. *J. Clin. Oncol.* **19**, 1728–1733 (2001).
 64. Institute of Medicine. *Transforming Clinical Research in the United States: Challenges and Opportunities: Workshop Summary* (The National Academies Press, 2010).
- Part of a report from a workshop at which issues relating to clinical-trial-recruitment statistics were presented and specific challenges were identified.**
65. Cancer Research UK. Stratified medicine and the lung cancer ‘Matrix’ trial — part of a cancer care revolution. *Cancer Research UK* <http://scienceblog.cancerresearchuk.org/2014/04/17/stratified-medicine-and-the-lung-cancer-matrix-trial-part-of-a-cancer-care-revolution> (2014).
 66. Tam, A. L. *et al.* Feasibility of image-guided transthoracic core-needle biopsy in the BATTLE lung trial. *J. Thorac. Oncol.* **8**, 436–442 (2013).
 67. Seguin, L. *et al.* An integrin $\beta 3$ –KRAS–RafB complex drives tumour stemness and resistance to *EGFR* inhibition. *Nature Cell Biol.* **16**, 457–468 (2014).
 68. I-SPY-2 Clinical Trials. About. *I-SPY 2 TRIAL* <http://ispy2.org/about> (2015).
 69. Barker, A. D. *et al.* I-SPY 2: an adaptive breast cancer trial design in the setting of neoadjuvant chemotherapy. *Clin. Pharmacol. Ther.* **86**, 97–100 (2009).
 70. National Institutes of Health. Molecular profiling-based assignment of cancer Therapy for patients with advanced solid tumors. *National Institutes of Health Clinical Center* http://clinicalstudies.info.nih.gov/cgi/detail.cgi?A_2013-C-0105.html.
 71. TrialReach. Clinical study for patients with cancer (Ve-Basket 120326). *TrialReach* <http://trialreach.com/study/clinical-study-for-patients-with-cancer-ve-basket-CT120326/> (2012).
 72. Worldwide International Networking. WIN Clinical Trials/Scientific Projects. *Worldwide International Networking in personalised cancer medicine* <http://www.winconsortium.org/page.jsp?id=104>.
 73. André, F. *et al.* Comparative genomic hybridisation array and DNA sequencing to direct treatment of metastatic breast cancer: a multicentre, prospective trial (SAFIR01/UNICANCER). *Lancet Oncol.* **15**, 267–274 (2014).
 74. Medical Research Council Clinical Trials Unit. Welcome to FOCUS4. *FOCUS4 Molecular selection of therapy in metastatic colorectal cancer: a molecularly stratified randomised controlled trial programme* <http://www.focus4trial.org/> (2014).
 75. Biankin, A. V. & Hudson, T. J. Somatic variation and cancer: therapies lost in the mix. *Hum. Genet.* **130**, 79–91 (2011).
- A review article that addresses the challenges presented by the molecular diversity in cancer that is uncovered through genomic sequencing.**
76. Eisenhauer, E. A. *et al.* New response evaluation criteria in solid tumours: revised RECIST guideline (version 1.1). *Eur. J. Cancer* **45**, 228–247 (2009).

Acknowledgements The authors would like to thank A. Ewing for her assistance in compiling the manuscript. They also thank L. Musgrove, D. Chang and P. Bailey for proofreading the manuscript and for their helpful suggestions.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of this paper at go.nature.com/ultylj. Correspondence should be addressed to A.V.B. (andrew.biankin@glasgow.ac.uk).

Progress and challenges in probing the human brain

Russell A. Poldrack¹ & Martha J. Farah²

Perhaps one of the greatest scientific challenges is to understand the human brain. Here we review current methods in human neuroscience, highlighting the ways that they have been used to study the neural bases of the human mind. We begin with a consideration of different levels of description relevant to human neuroscience, from molecules to large-scale networks, and then review the methods that probe these levels and the ability of these methods to test hypotheses about causal mechanisms. Functional MRI is considered in particular detail, as it has been responsible for much of the recent growth of human neuroscience research. We briefly review its inferential strengths and weaknesses and present examples of new analytic approaches that allow inferences beyond simple localization of psychological processes. Finally, we review the prospects for real-world applications and new scientific challenges for human neuroscience.

The way that we conceptualize brain function has always been constrained by the methods available to study it. Studies of patients with focal brain lesions in the nineteenth century led to the view of the brain as a collection of focal centres specialized for particular cognitive abilities, such as ‘Broca’s area’ for speech production. The development of neurophysiological recording techniques in the twentieth century led to Barlow’s ‘neuron doctrine’, according to which the functions of individual neurons can be extrapolated to explain the function of the brain as a whole. The cognitive neuroimaging studies of the 1980s focused on subtractive comparisons between cognitive tasks meant to isolate specific cognitive operations, and led to a relatively modular view of brain function as involving localized and separable regions that implement elementary mental operations.

The methods of contemporary human neuroscience have provided a much more complex and nuanced view of the human brain as a dynamic network with multiple levels of organization, in which function is characterized by a balance of regional specialization and network integration. Although current methods are limited in their utility for studying brain function at fine-grained levels of organization (such as single neurons or cortical columns), human neuroscience has nonetheless made remarkable progress in understanding basic aspects of functional organization, and with this have come a number of applications to address real-world problems. Our goal here is to review the current state of human neuroscience, focusing on what kinds of questions can and cannot be answered using current techniques and how those answers are relevant to real-world applications.

How can we study the human brain?

Methods for studying human brain function can be organized according to the kinds of mechanistic insights that each technique provides. As shown in Table 1 the first characteristic is the level of mechanism captured by the method. Mechanisms range from the molecular level (neurotransmitters and receptors) to large-scale networks (the dynamic integration and coordination of different functional areas of the brain). Although this distinction is related to physical scale, it does not depend on the method’s spatial resolution *per se*. For example, positron emission tomography (PET) using neurotransmitter ligands measures molecular mechanisms, even though its spatial resolution is on the order

of one centimetre. The second characteristic is the ability of each method to elucidate the mechanistic role of an observed brain molecule, cell, region or network in a mental function of interest. By mechanism we mean the causal chain of events that result in the realization of a function. To fully understand human brain function is to know the causal chains of events at the molecular, cellular, population, and network levels that give rise to psychological function. For this reason, the power to identify causal relationships is a crucial dimension of difference among methods.

Some methods used in the study of human brain function provide relatively little insight into causal mechanisms. This includes methods that exploit naturally occurring variation by observing the strength of association between individual differences in brain function and behaviour. Analysis of relationships between behavioural traits, genes, brain structure, and brain function exemplify this approach (see Box 1 for a discussion of genomic approaches). For many important psychological phenomena, from effects of life history to personality traits, we are limited to observational methods. For example, individual differences in the personality trait of impulsiveness have been associated with differences in striatal dopamine release¹, functional MRI (fMRI) activation², and cortical grey matter volume³. Observed associations between neural and psychological traits do not necessarily imply a causal relationship, as these associations could result from an unmeasured third variable that independently influences the two measures. Nevertheless, such associations provide a valuable starting point for theorizing about the neural mechanisms of human psychology, and their evidentiary value can be strengthened by measuring possible confounds to rule them in or out.

Although functional neuroimaging, electroencephalography/magnetoencephalography (EEG/MEG) and single-cell recordings are sometimes criticized as being purely correlative and therefore uninformative about mechanism, that criticism is only partly accurate. When psychological processes are experimentally manipulated by presenting a certain kind of stimulus and/or engaging the subject in a task, we can infer that any reliably elicited brain activity was caused by performing these psychological functions. We cannot, however, infer with confidence that the observed brain activity is causally responsible for the psychological process under study. Despite this limitation (which is shared by neuronal recordings in non-human animals), neuroimaging studies

¹Department of Psychology, Stanford University, Stanford, California 94305, USA. ²Center for Neuroscience & Society, University of Pennsylvania, Philadelphia, Pennsylvania 19104, USA.

Table 1 | An overview of the levels of analysis and levels of causal inference afforded by different human neuroscience methods

		Level of mechanism			
		Molecules	Cells	Populations	Networks
Strength of causal evidence	Purely observational (associations do not necessarily imply causal relations between mind and brain)	Genetic associations with behaviour, brain function or brain structure Postmortem studies of gene expression Correlations of MRI spectroscopy or PET ligand imaging with psychological traits		Structural morphometry correlated with psychological traits	Resting functional connectivity (fMRI, EEG/MEG) or structural connectivity (sMRI, DTI) correlated with psychological traits
	Manipulate psychological process and observe brain (neural measures may be epiphenomenal)	Task modulation studies using PET with neurotransmitter ligands or MRI spectroscopy	Intracerebral recording in surgical patients	Task activation studies (PET, fMRI, EEG/MEG) Representational analysis (fMRI, EEG/MEG) Computational neuroimaging (fMRI, EEG/MEG)	Task-based functional connectivity (fMRI, EEG/MEG)
	Manipulate brain and observe psychological results (demonstrates causal effect of neural system in behaviour)	Pharmacological manipulation (including hormones)	Direct brain stimulation in surgical patients	Focal cortical lesions Transcranial magnetic stimulation Transcranial electrical stimulation Cortical surface electrode stimulation in surgical patients	Disconnection/white matter lesions

DTI, diffusion tensor imaging; EEG/MEG, electroencephalography/magnetoencephalography; fMRI, functional MRI; MRI, magnetic resonance imaging; PET, positron emission tomography; sMRI, structural MRI.

in which psychological processes are manipulated comprise the majority of current human neuroscience research, and have advanced our understanding of human brain function, as we will discuss in more detail below.

BOX 1

Challenges of merging neuroimaging and genomics

The substantial heritability of many psychological functions has driven great interest in finding genetic underpinnings of individual differences in neural function. Twin and family studies have demonstrated significant heritability for both task-related BOLD responses⁹¹ and resting-state functional connectivity⁹² in fMRI. In the past decade, a large number of studies have also reported associations between BOLD signals and common variants in candidate genes. Unfortunately, this approach has generally been unsuccessful in identifying genetic associations that are replicated in genome-wide association studies (GWAS). For example, a striking finding from the first well-powered GWAS of genetic variants associated with brain volume was that none of the associations previously identified through candidate gene studies were replicated at the genome-wide level⁸⁸. Similarly, candidate gene associations with cognitive function (such as the association between polymorphisms in the *COMT* gene and working memory) and brain activation have generally not been confirmed in meta-analyses, and are subject to a substantial degree of publication bias^{93,94}. Like for many other areas of genetics, this suggests that genome-wide approaches are the most likely to lead to reliable identification of common variants related to brain structure and function. However, GWAS approaches require large samples (in the tens of thousands) which are very difficult to amass for task-based fMRI studies; for that reason, GWAS-based approaches to probing the human brain will likely be limited to resting-state fMRI and structural MRI. Other strategies, such as targeted studies investigating rare variants of large effect identified using genome sequencing or studies using gene expression in peripheral tissues may have greater utility for genetic studies of task-based fMRI. Task-based fMRI may also be used to further investigate candidate variants identified on the basis of GWAS studies of psychiatric disorders or population variability.

More decisive evidence concerning causal necessity can be obtained by manipulating the brain itself to assess the resulting effect on the psychological process in question. Naturally occurring or surgical lesions, which provided the basis for most of what we knew about human brain function before the advent of neuroimaging, are still of great interest because they provide insight into the causal necessity of specific brain regions or connections. More recently developed methods of brain stimulation allow for reversible inhibition or excitation of a brain area, thereby expanding our ability to test the causal role of brain regions in the mechanisms of human thought and action. Deep brain stimulation (DBS) provides the most precise method for targeted stimulation by using surgically implanted electrodes, but is limited to situations where patients are undergoing implantation for medical reasons. Use of non-invasive brain stimulation for research purposes has grown rapidly in recent decades, starting with transcranial magnetic stimulation (TMS), in which pulsed magnetic fields induce currents in the brain. Various forms of transcranial electric stimulation (TES), in which current is delivered using external electrodes, have also been used, of which the most common variant is transcranial direct current stimulation (tDCS). Unlike DBS, non-invasive brain stimulation generally affects larger and more superficial areas of the brain, but researchers are seeking to improve spatial resolution with new magnetic coil shapes for TMS and new electrode configurations for tDCS. Focused ultrasound is also being explored as a means to stimulate more precisely delimited brain regions⁴. Pharmacological agonists and antagonists of particular neurotransmitter systems can be used to experimentally manipulate the human brain at the molecular level, although with imperfect specificity⁵. By combining each of these manipulations of brain function with functional brain imaging, one can leverage the causal information obtained through pharmacological challenges or brain stimulation. For example, the causal role of activity in specific brain regions, identified using fMRI, for a particular function has been tested by brain stimulation, using both direct cortical stimulation (for example, ref. 6) and TMS⁷.

New capabilities of fMRI

Because fMRI has become the main method for the study of human brain function, our review focuses on this method and new ways of using it. In the last two decades, fMRI has transitioned from a newly developed technique for revealing neuronal activity to being the workhorse method of cognitive neuroscience (see the recent special issue of *Neuroimage* on

the first 20 years of fMRI⁸). Much has been learned about the biological mechanisms underlying blood oxygen level dependent (BOLD) signals^{9,10}, but still much remains to be understood, such as the roles of specific glial and neuronal cell types in the coupling of neuronal activity to blood flow (for example, refs 11, 12). This limited physiological understanding poses problems for the interpretation of fMRI data. In particular, although fMRI signals often correlate strongly with both action potentials ('spikes') and local field potentials, they are largely reflective of post-synaptic processes, and in some cases they can be dissociated from spiking altogether¹³. The relative sensitivity of fMRI to post-synaptic processes as opposed to spiking has been seen as a drawback by some who view spikes as the essence of brain function, but it is worth noting that this discovery has actually rekindled interest in the analysis of local field potentials in electrophysiology (where these signals have long been discarded) (for example, ref. 14), and suggests that fMRI may sometimes be sensitive to subthreshold signals that would be missed by analysis of spikes only. Uncertainties in relating fMRI to psychological, as well as physiological, processes have also been debated, and progress has been made on this front too. From experimental approaches such as adaptation paradigms for probing representations to analyses of functional connectivity, fMRI is routinely used to answer questions about mind–brain relationships that go far beyond localization¹⁵. Here we discuss three examples of new approaches to understanding human brain function with fMRI that address questions of representation, computational processes and network interactions across the brain.

Representational analyses

Early work in neuroimaging focused largely on 'brain mapping'—identifying regions based on the mental processes that cause them to be activated. This approach has provided a large body of reliable associations between function and structure, but has not been particularly successful in providing new insights into how psychological functions are implemented¹⁶. However, two relatively recent approaches, known as multi-voxel pattern analysis (MVPA)¹⁷ and representational similarity analysis (RSA)¹⁸, can more directly relate psychological contents to brain function (Fig. 1). MVPA involves the use of methods from the field of machine learning to decode or predict psychological states from patterns of brain activation across voxels (hence the term 'brain-reading'). Since its introduction more than a decade ago, MVPA has been used in a number of domains to demonstrate the predictive ability of fMRI activation patterns. Perhaps the most impressive are demonstrations of the ability to successfully reconstruct visual scenes¹⁹ and faces²⁰ from BOLD activity patterns; similar advances have been made for higher cognitive functions such as word meaning²¹. These studies go beyond simply differentiating between experimental conditions, as they show how the underlying representational spaces relate to brain activity; for example, using a related approach known as voxel-wise modelling, Huth and colleagues²² developed a model that estimated the response at each location on the cortical surface to a large number of visual and semantic features present in natural movies (Fig. 2). MVPA approaches have also provided new insights into the neural organization of cognitive functions. For example, MVPA has informed our understanding of the mechanisms of visual attention, by showing that attention changes both the representation of stimuli across regions of visual cortex as well as the mutual information between regions²³. In the domain of memory, MVPA has been used to show that competition between memory representations in working memory leads to poorer subsequent memory for those items, demonstrating a nonmonotonic relationship between competition and subsequent memory²⁴.

Whereas MVPA is generally used to decode individual psychological states, RSA instead asks how the patterns of brain activity evoked by different stimuli are related to one another, and thus provides the means to directly address questions of how mental representations are implemented in the brain. RSA has enabled the demonstration of

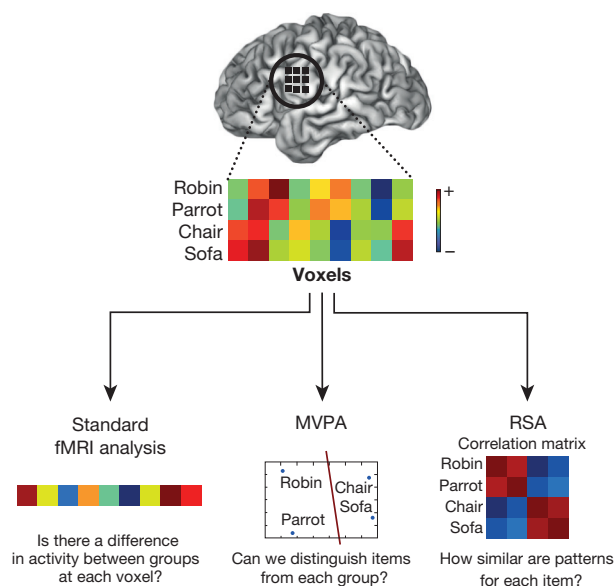


Figure 1 | Different approaches to the analysis of fMRI data. This example depicts data from a hypothetical study in which four different stimuli were presented (two birds and two items of furniture) and response measured for each item across nine voxels; intensity of activity is depicted from blue (negative) to red (positive). The standard univariate fMRI analysis approach would examine the difference at each voxel between the averages of the two categories. Multi-voxel pattern analysis (MVPA) examines the multidimensional relationship between patterns of activity, in this case projecting the nine-dimensional space of voxel patterns (the voxel vector) into a two-dimensional space and identifying a boundary that separates items from the two classes. Representational similarity analysis (RSA) examines the correlations between activity patterns for each item, in this case showing that items within category show a high correlation (red), whereas the correlation of items between categories is low (blue).

direct isomorphisms between psychological representations of stimuli (such as the similarity or typicality of objects) and the neural patterns associated with those stimuli^{25,26}. Because psychological theories often make predictions regarding the similarity of different stimuli, RSA has also enabled the direct testing of theories, such as theories about how categories are represented²⁷ and theories of how repeated experiences lead to enhanced learning²⁸. RSA can be applied to any kind of multi-dimensional data, and this has enabled the demonstration of systematic mappings of visual object representations between humans (using fMRI) and non-human primates (using electrophysiological recordings)²⁹—an example that highlights how human neuroscience can also help to establish more direct parallels with findings in non-human models, allowing insights to filter in both directions.

Although much MVPA and RSA work (as depicted in Fig. 1) has focused on the representations found in localized brain regions, these methods are equally useful for assessing representations that are spread across the brain. For example, recent work has shown that mental states such as physical pain can be decoded by analysis of patterns of activation across brain regions³⁰.

The legitimate enthusiasm about these methods is tempered by lingering questions regarding the interpretation of multivariate analyses^{31,32}. In addition, recent work combining electrophysiology and fMRI in non-human primates has demonstrated that the sensitivity of MVPA is limited by the spatial characteristics of the neuronal representations that code for particular features, such that some kinds of neuronal patterns may be more difficult to decode using MVPA than others³³. Finally, it is important to stress that, like standard neuroimaging approaches, MVPA and RSA approaches do not inform about causal mechanisms.

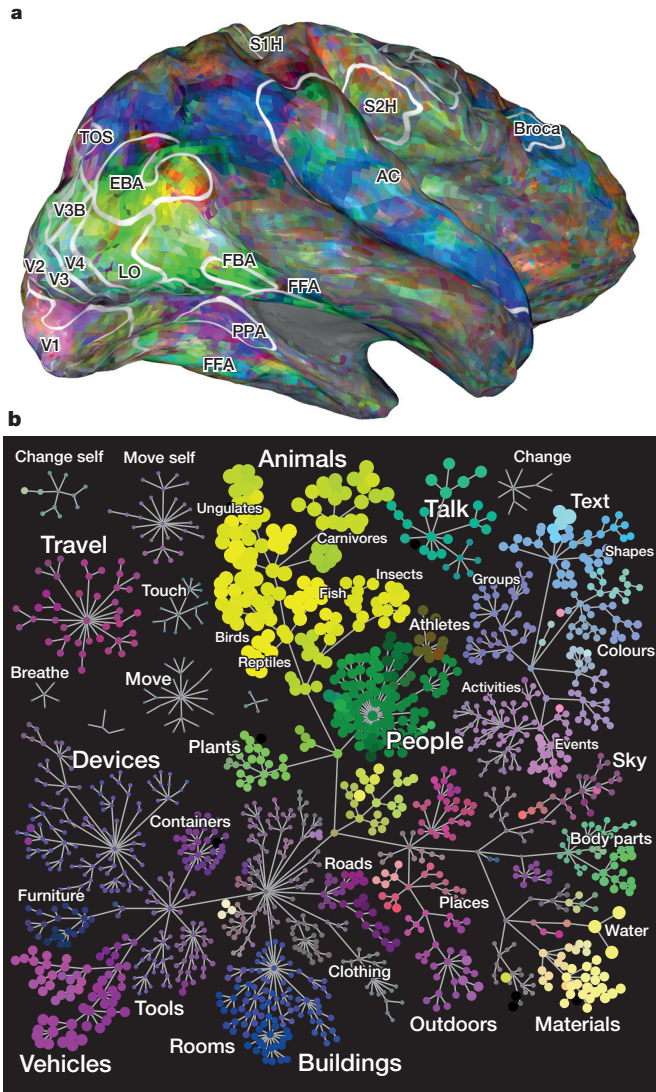


Figure 2 | A mapping of high-dimensional semantic space onto the cortical surface. Here, voxel patterns for 1,705 different action and object categories, based on brain activity obtained during viewing of natural movies²² are mapped onto the cortical surface image generated using online browser at (<http://gallantlab.org/semanticmovies/>). **a**, Mapping of semantic categories to each point on the surface; the colours on the surface map correspond to the semantic map in panel **b**. **b**, A depiction of the semantic space derived from all semantically selective voxels. Categories that have similar colours in the semantic space are represented in similar patterns of voxels in the brain. Data from ref. 22.

Integrating fMRI and computational modelling

Computational models play a central role in our understanding of both cognitive and brain functions and, increasingly, of the relationship between the two. By making assumptions explicit, computational models enable more direct testing of theories, as well as providing the means to link computations at the neuronal level with higher-order functions. An example of an area in which substantial progress has been made using this approach is reinforcement learning, in which an animal selects actions and learns from the rewards gained from those actions. Computational models of reinforcement learning (RL) have long played a central role in artificial intelligence and psychology, and the discovery by Schultz and colleagues³⁴ that dopamine neurons appear to signal one of the important quantities in these models (reward prediction error) has brought these models to the forefront of the neuroscience of decision making. For example, a set of publications in 2003 applied RL models to neuroimaging data and thereby identified correlates of reward predic-

tion error signals in dopaminergic target regions such as the ventral striatum^{35,36}. Subsequent neuroimaging work has established that there are multiple RL signals in the brain, some reflecting the simple association between actions and values (known as 'model-free' RL) and others reflecting more complex contextual and hierarchical learning processes (known as 'model-based' RL)^{37,38}. Similarly, in the study of memory, progress has been made in the mapping of medial temporal lobe subregions to specific computational operations such as pattern completion and pattern separation (for example, ref. 39). In each of these domains, the computational interpretation of neuroimaging signals has been greatly enhanced by parallel studies in non-human animals, allowing imaging signals to be linked more directly to direct measures of neuronal activity.

Functional connectivity analysis and resting-state fMRI

Perhaps the most revolutionary development to arise from human neuroimaging research is the realization that the resting brain is far from quiescent, and that important insights into brain function can be gained from studying the correlated fluctuations of signals across the brain at rest. Much of the research into the resting state has focused on a set of regions (including anterior and posterior midline regions, lateral temporoparietal cortex, and the medial temporal lobe, known as the 'default mode' network⁴⁰) that are consistently less active during performance of difficult tasks⁴¹, and are functionally connected in the resting state⁴². Similar patterns of resting connectivity have been observed in non-human primates⁴³ and awake rodents⁴⁴, suggesting that they reflect fundamental principles of mammalian brain organization. There is also growing evidence that these networks may be important in brain disorders. For example, the posterior portion of the default mode network appears to play a critical role in the memory deficits observed in Alzheimer disease, showing a convergence of amyloid deposition, structural atrophy, and decreased metabolic activity⁴⁵.

Data collected in the resting state can provide insights into the broader functional organization of the brain as well. In particular, the organization of resting state signals bears a close relation to the organization of brain activity evoked by mental tasks. For example, Smith *et al.*⁴⁶ used independent component analysis to identify spatially independent sets of voxels from resting-state fMRI data and from task-based data (obtained from the Brainmap meta-analytic database), and demonstrated that the components extracted from resting-state fMRI showed a high degree of concordance with those extracted from task-based data. The overlap between resting-state and task-based functional organization can also be seen within individuals; for example, the longitudinal examination of a single individual revealed reliable spatial parcellation of activity in the cerebral cortex (using resting fMRI data) that mapped systematically to the activation patterns observed across a large number of task measurements⁴⁷.

Despite the substantial excitement around resting-state fMRI findings, numerous concerns have been raised about their interpretation. In particular, there are lingering questions regarding the ways in which artefacts related to head motion and physiological fluctuations may influence estimates of resting state connectivity, and whether common data analytic methods may induce systematic artefacts^{48,49}. In addition, potential confounds such as light sleep⁵⁰ may drive differences in resting state signals. The unconstrained nature of resting-state fMRI is a double-edged sword; it is potentially very useful for the study of clinical groups for whom task performance may be difficult, but at the same time, it is not possible to determine whether group differences reflect fundamental differences in functional connectivity or relative differences in the ongoing mental content of different groups during rest (see ref. 51).

Applications of human neuroscience

With the development of new methods have come attempts to apply them to real-world problems, in both medical and non-medical contexts. (See Box 2 for a discussion of the ethical, legal, and societal issues raised by these applications.)

BOX 2

Ethical, legal and societal impact of human neuroscience

As the methods of human neuroscience find broader application, they affect human life in new ways. The field of 'neuroethics' is concerned with ethical, legal and societal issues raised by these new applications⁹⁵.

Two kinds of problems have emerged from the increasing ability of brain imaging to reveal aspects of individual psychology: problems that arise from the current and imminent capabilities of these methods, and problems that arise from their lack of claimed capabilities. To the extent that imaging can predict important personal characteristics such as health status, academic achievement, and criminal behaviour, its use must be managed with care to protect privacy and avoid discrimination⁹⁶. To the extent that imaging cannot provide help with high-stake problems, the public should be protected from claims that it can. For example, a seemingly 'scientific' method for detecting lies or diagnosing psychiatric disorders^{66,97} has a strong appeal to the general public who cannot be expected to appreciate the gap between what is claimed and what is established fact.

New ways of changing brain function pharmaceutically, and with electromagnetic stimulation, also raise new ethical issues. Of course, humanity has long manipulated brain function to modify mental states using substances such as alcohol and caffeine. However, psychopharmacology has broadly penetrated our everyday lives and the scope of psychiatric diagnoses and treatment has expanded—a societal shift that some find troubling⁹⁸. Furthermore, many now use psychoactive drugs purely for enhancement of healthy brain function rather than to treat a medical condition⁹⁹. Aside from issues of safety and efficacy, brain enhancement raises issues of fairness (is it akin to doping in sports?), justice (will the ability to access enhancements widen the already existing gaps between haves and have-nots?) and social standards (will unenhanced job performance become sub-standard?).

Non-invasive brain stimulation is the newest method for brain enhancement. Simple transcranial electrical stimulation (for example, tDCS) devices are available to consumers at relatively low cost and regulation is minimal¹⁰⁰. Given the public interest in this method and the rudimentary state of knowledge about its effects, it is crucial that the safety and efficacy of these methods are established. The efficacy of cognitive enhancement with tDCS is hotly debated¹⁰¹ and whether long-term use of tDCS is safe has yet to be studied. In addition, neuroethical issues of fairness, justice and social standards mentioned above also apply to enhancement of brain function by brain stimulation.

Brain disorders

The methods of human neuroscience hold particular promise for understanding and treating psychiatric disorders, because these disorders do not have clear analogues in non-human animals, and animal models currently used for preclinical screening of potential therapies are increasingly regarded as being inadequate⁵². In the absence of valid animal models, it becomes all the more crucial to apply new methods for understanding human brain function and dysfunction. The goal of improving the treatment of neuropsychiatric disorders is made even more challenging because of our current diagnostic system. Although depression, schizophrenia, autism and other serious psychiatric disorders have long been considered disorders of the brain, they are still diagnosed exclusively by behavioural signs and symptoms. These diagnostic criteria do not seem to have clear relations to the biological processes that would be targeted by new medical treatments.

In response to this problem, an alternative way of systematizing psychiatric disorders has been developed—the NIMH Research Domain

Criteria (RDoC)⁵³—that describes disorders according to impairments in specific functional systems of the brain (such as fear or reward learning) and at different levels of mechanism of the kind represented in Table 1 (for example, molecules or circuits). RDoC characterizations cut across traditional diagnostic categories and are intended to capture the underlying pathophysiology more accurately. Given the multiple levels of mechanism captured by the RDoC, the system encourages research with a broad array of methods to identify potentially targetable dysfunctions.

The application of several human neuroscience methods has led to the development of targeted treatments, for example, in the field of depression. Functional imaging studies have highlighted the role of the subgenual anterior cingulate cortex in a network of regions involved in mood, leading Mayberg and colleagues to use deep brain stimulation in this area to regulate mood in depressed patients⁵⁴. Lateral prefrontal regions, implicated through imaging studies in depression, have been targeted with non-invasive brain stimulation, including the FDA-approved use of TMS for treatment-resistant depression. Functional neuroimaging can itself be used as a treatment, by providing patients with a real-time measure of regional brain activity to use as a biofeedback signal. This approach is being tested for the treatment of chronic pain, depression and addiction⁵⁵. In contrast, neuroimaging has not so far been very successful in aiding differential diagnosis of disorders in terms of current diagnostic categories. A recent large meta-analysis identified a set of regions in which structural abnormalities were consistently associated with psychiatric disorders, but found very little specificity for individual disorders⁵⁶, consistent with the notion that current diagnostic distinctions are not biologically realistic categories.

Another approach to the discovery of therapeutic targets is the use of genetic association studies to identify sets of genes that are associated with a disorder and that together may indicate particular molecular pathways underlying the disorder. Although the numbers of subjects needed to establish reliable genetic associations is daunting, progress has been made through large international collaborations. For example, Psychiatric Genomics Consortium has to date identified more than 100 common genetic variants reliably implicated in schizophrenia⁵⁷. Imaging can also be used to develop endophenotypes (or intermediate phenotypes) that may bear a closer relation to the effect of a gene variant than does disease diagnosis, as well as to mitigate the problem of heterogeneity within conventional diagnostic categories (see Box 1).

It may be less surprising that the methods developed for human neuroscience research have been applied in the diagnosis and treatment of neurological diseases, but at least two recent developments deserve mention here. Studies of Alzheimer disease at mechanistic levels from molecules to systems have improved diagnostic accuracy and have enabled a degree of prediction before clinical signs of the disease⁵⁸. Molecular biomarkers from blood and CSF, and patterns of brain activity and structure have revolutionized clinical research in this area by facilitating trials of preventive treatment and by providing intermediate phenotypes as early gauges of effectiveness. Disorders of consciousness following severe brain damage are another area of clinical neuroscience for which neuroimaging shows promise. Some patients who have been diagnosed as being in the vegetative state can follow commands to imagine actions that activate specific areas of the brain in much the same way as healthy control subjects do, and can even use these imagined actions to answer questions (for example, "Do you have any brothers? If yes, imagine playing tennis, if no, imagine walking through your house.")⁵⁹. Thus, neuroimaging offers new insights into the assessment of consciousness, as well as the distinct problem of prognosis, in severely brain-damaged patients.

Predicting behaviour

The ability to predict future behaviour is of value in almost every sphere of human activity. Although it has often been said that "the best predictor of future behaviour is past behaviour," in some cases brain imaging can improve our ability to predict future behaviour, over and

above what we can do with behavioural history. Marketing professionals were among the first to attempt to predict behaviour using brain imaging. Recognizing the limitations of focus groups and other traditional methods to discern what consumers want, they have used functional neuroimaging to predict the effects of different advertising campaigns, packaging, and other factors on consumer behaviour, based on the premise that activity in the brain's reward or motivation centres may be a more direct measurement of wanting than are verbal self-reports⁶⁰. Although most of this work is conducted by and for corporations aiming to improve sales rather than share scientific knowledge, published academic studies have begun to lend some credence to the potential of neuromarketing. For example, when teenage subjects were scanned while listening to unfamiliar songs, the reward system activity evoked by the songs, but not the subjects' ratings of their likeability, was predictive of sales of the songs over the subsequent three years⁶¹.

Prediction is also important outside of business. Falk and colleagues have adapted neuromarketing methods for the purpose of creating more effective public service announcements. They showed that brain responses (but not ratings) to an anti-smoking advertisement were predictive of subsequent call volume to an anti-smoking hotline⁶². Gabrieli *et al.*⁶³ recently summarized evidence concerning neuroimaging-based prediction in domains ranging from healthful eating to criminal recidivism, including numerous examples of prediction of educational outcomes. Indeed, neuroimaging can predict future academic skills over and above traditional behavioural predictors, thus enabling earlier and more appropriate interventions to address individual children's reading and math difficulties. These authors also pointed out a number of methodological challenges in neuroimaging-based prediction of behaviour, including the need to develop and test predictions with different samples, to avoid the 'optimism bias' that occurs when predictions are tested in the same population from which they were generated.

Human neuroscience in the courtroom

In recent years the methods of human neuroscience have found their way into the courtroom. Perhaps the most obvious, but also the most misunderstood, role for neuroscience is in helping to determine criminal responsibility. Proving that a criminal act may have had a neural cause is not in itself exculpatory, as every human act is caused by the brain⁶⁴. However, to the extent that neuroscience can provide evidence of mental dysfunction (for example, a tumour in the frontal cortex that may have impaired the ability to control behaviour), immaturity or other psychological grounds for reduced criminal responsibility, it is potentially relevant and has been used. For example, the Supreme Court explicitly cited neuroscience evidence in its decision in *Graham v. Florida* to abolish life in prison without parole for juveniles who commit non-homicidal offences. It is more difficult to make legal arguments for applying neuroimaging evidence to individual cases because most findings from neuroimaging research are generalizations based on groups of people and may therefore not allow reliable inferences regarding individuals⁶⁵. Nevertheless, neuroimaging scans from defendants are sometimes presented in the sentencing phase of criminal trials as grounds for mitigation of the sentence, as weaker evidentiary standards apply in the sentencing phase.

Neuroimaging can be applied in ways other than determining degree of responsibility. Lie detection is one example that has been pursued in legal contexts, although it has not so far been admitted into US courts and has yet to demonstrate validity, reliability or resistance to countermeasures outside of the laboratory⁶⁶. Another application concerns pain: brain-based biomarkers for pain would help discriminate real suffering from malingering—a pivotal issue in many lawsuits—and have been admitted as evidence in at least one US case⁶⁷.

Challenges and future directions for neuroimaging

The field of neuroimaging is growing rapidly, and there are a number of exciting new directions on the horizon.

New technologies for imaging and manipulating the human brain

Rapid advances in non-human neuroscience have been driven by the development of technologies that measure and manipulate brain function with increasing precision. Human neuroscience has lagged in this respect, in part because of the ethical challenges associated with direct manipulation and neuronal recording of the human brain. However, in response to the urgent need for better treatments for psychiatric disorders, research is underway with the aim to design implantable systems for sensing and modulating human brain networks⁶⁸. The development of optogenetic and 'opto-fMRI' approaches in non-human primates⁶⁹ suggests that these methods may one day become feasible for use in human studies, and it is likely that electrical brain stimulation will eventually be supplemented with optogenetic approaches. Although such invasive techniques will likely only be used in rare clinical cases (that is, patients are undergoing implantation for medical reasons), they have the potential to provide much greater specificity in circuit mapping.

fMRI will probably remain the principal neuroimaging method in humans in the foreseeable future. However, the ongoing BRAIN initiative in the United States⁷⁰ is providing substantial funding to develop entirely new techniques for imaging of brain function, and a significant proportion of this funding will go specifically towards the development of new methods for imaging the human brain. In addition, new developments in MRI have greatly increased the utility of standard MRI systems. For example, multiband imaging techniques⁷¹ have enabled a several-fold increase in the temporal resolution of fMRI acquisitions, and higher MRI field strengths (7 tesla and higher) hold promise to enable improvements in spatial resolution as well (for example, ref. 72). There is thus great reason to be optimistic that methodological limits will continue to be pushed in the future.

Additional insight into human brain function will likely come from the study of postmortem human brains, which has long been a staple method for the characterization of anatomical structure and study of brain disorders. New techniques have enhanced the ability to visualize the structure of human brain tissue (Fig. 3). For example, optical coherence tomography has been used to image *ex vivo* human cortical tissue, providing high-resolution imaging of cytoarchitecture with less distortion than standard microscopy techniques⁷³. The first whole-brain atlas of genome-wide gene expression in postmortem human brains⁷⁴ has provided an important resource for understanding how gene expression relates to brain function; for example, the maps from this project have

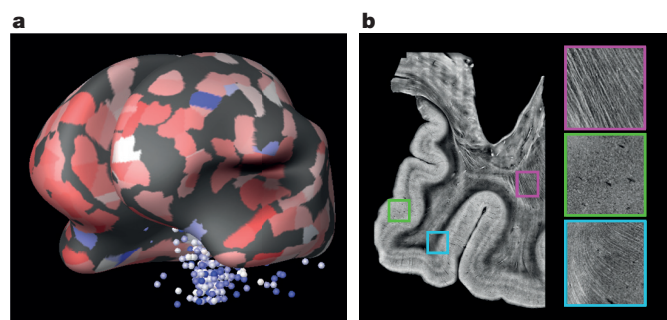


Figure 3 | New methods for characterizing the postmortem human brain.

a, A map of expression of the serotonin receptor 3B displayed on the reconstructed cortical surface in one individual from the Allen Brain Atlas Human Brain data set (generated using data from <http://human.brain-map.org/>). **b**, Optical coherence tomography imaging of the human brain (2.9 in-plane resolution). Large panel presents an average intensity projection in depth over 300; inset zooms are maximum intensity projections over 300, showing fibres in the white matter (pink inset), fibres arcing through the subcortical junction to insert into the cortex (cyan inset), and neurons in the cortex (bright spots in the green inset). Image courtesy of Bruce Fischl, Caroline Magnain and David Boas, Massachusetts General Hospital.

been used to identify expression differences across different resting-state networks⁷⁵. Continued development of such resources will be essential for progress in understanding the genetic architecture of brain function and their relation to mental health disorders.

Connectomics

The Human Connectome Project⁷⁶ is nearing completion, and has already provided a rich database for the modelling of functional and anatomical connectivity of the human brain. However, fundamental challenges remain. For example, diffusion MRI provides the means to track white matter pathways (Fig. 4) and has been used to identify white

matter connectivity disruptions associated with cognitive disorders such as dyslexia⁷⁷; however, diffusion imaging has inherent biases that limit its ability to accurately track connections across the entire brain^{78,79}. The last decade has seen a proliferation of approaches to model functional connectivity on the basis of functional MRI data, though the dust has yet to settle regarding which methods are most effective (for example, ref. 80). To determine this, the analysis methods must be validated, which is challenging to do in humans but may be achieved using direct measurements of functional connectivity from invasive human approaches and non-human animals to validate the neuroimaging results. There is increasing evidence that at least in non-human primates

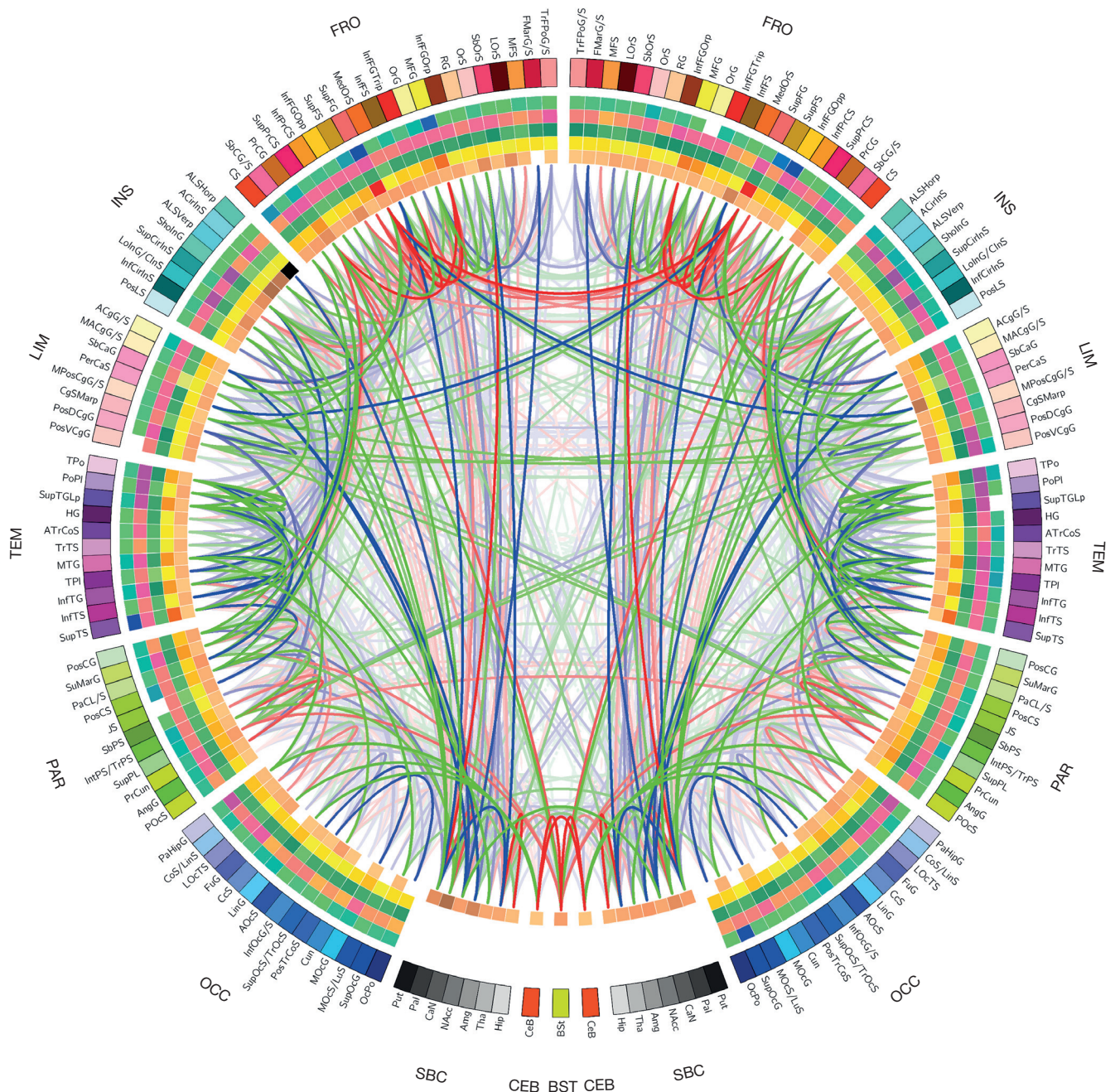


Figure 4 | A 'connectogram'⁹⁰ for an example healthy adult female subject. The outermost ring shows the various brain regions arranged by lobe (fr, frontal; ins, insula; lim, limbic; tem, temporal; par, parietal; occ, occipital; nc, non-cortical; bs, brain stem; CeB, cerebellum) and further ordered anterior (top) to posterior (bottom). The colour map of each region is lobe-specific and maps to the colour of each regional parcellation as determined using FreeSurfer.

The set of five rings (from the outside inward) reflect grey matter volume, area, thickness, curvature, and connectivity density. The lines inside of the circle represent the computed degrees of connectivity between segmented brain regions using diffusion tractography, with colour representing the relative fractional anisotropy of the connection (from blue to red). Image courtesy of Jack Van Horn, University of Southern California.

functional connectivity reflects anatomical connectivity as measured using either diffusion MRI⁸¹ or anatomical tract-tracing⁸²; but it remains an important challenge to establish the ways in which functional and diffusion connectivity measures converge or diverge.

Reproducibility of neuroimaging research

Large-scale meta-analyses have made it clear that neuroimaging results can be highly convergent across studies, to the degree that cognitive processes can be accurately inferred from individual subject data using decoders trained on meta-analytic data based on reported activation coordinates⁸³. However, the last few years have also seen increasing concern regarding the reproducibility of research findings in neuroscience, paralleling more general concerns about reproducibility of scientific results⁸⁴. These issues are particularly acute for neuroimaging given the high dimensionality of the data, relatively low statistical power of many studies⁸⁵, high degree of analytic flexibility in data analysis procedures⁸⁶, and potential for questionable research practices such as circular analysis procedures⁸⁷. The field of neuroimaging has been at the forefront of a number of developments that aim to improve reproducibility and the sharing of data are increasingly being embraced. The Alzheimer's Disease Neuroimaging Initiative (ADNI), International Neuroimaging Data Sharing Initiative (INDI), ENIGMA, and the Human Connectome Project together have shared thousands of neuroimaging data sets and this has enabled a number of novel discoveries. For example, data sharing by the ENIGMA consortium has enabled the first well-powered genome-wide association study of brain volume⁸⁸, identifying replicated associations between brain volume and several common genetic variants. In addition, nearly all of the main software packages for neuroimaging data analysis are free and open source, providing transparency and reproducibility in data analysis across groups, and the publication of fully reproducible analysis workflows has begun (for example, ref. 89). The increasing use of machine learning methods, with their focus on out-of-sample generalization rather than statistical significance, is also leading to a greater emphasis on achieving reproducibility.

Outlook

The use of new tools for imaging and manipulating the brain will continue to advance our understanding of how the human brain gives rise to thought and action. The combination of myriad methods with different and complementary strengths and weaknesses will allow neuroscientists to develop a multilevel understanding of the brain, spanning from molecules to large-scale networks. New analysis methods have advanced fMRI beyond 'blobology' and will provide direct insight into the mapping of mental and neural representations, while newer analysis and acquisition methods will offer other novel insights into the relation of mind and brain. fMRI and other human neuroscience methods will continue being applied to solve real-world problems, within medicine and beyond. Although some of these applications are currently premature relative to the demonstrated capabilities of the methods, it is clear that the new methods of human neuroscience will have much to offer science and society.

Received 31 May; accepted 4 September 2015.

1. Buckholz, J. W. *et al.* Dopaminergic network differences in human impulsivity. *Science* **329**, 532 (2010).
2. Plichta, M. M. & Scheres, A. Ventral-striatal responsiveness during reward anticipation in ADHD and its relation to trait impulsivity in the healthy population: a meta-analytic review of the fMRI literature. *Neurosci. Biobehav. Rev.* **38**, 125–134 (2014).
3. Schilling, C. *et al.* Common structural correlates of trait impulsiveness and perceptual reasoning in adolescence. *Hum. Brain Mapp.* **34**, 374–383 (2013).
4. Legon, W. *et al.* Transcranial focused ultrasound modulates the activity of primary somatosensory cortex in humans. *Nature Neurosci.* **17**, 322–329 (2014).
5. Chamberlain, S. R., Müller, U., Robbins, T. W. & Sahakian, B. J. Neuropharmacological modulation of cognition. *Curr. Opin. Neurol.* **19**, 607–612 (2006).
6. Parvizi, J. *et al.* Electrical stimulation of human fusiform face-selective regions distorts face perception. *J. Neurosci.* **32**, 14915–14920 (2012).

This study uses the combination of fMRI and intracranial electrical stimulation to demonstrate the causal role of fusiform regions in face perception.

7. Chen, A. C. *et al.* Causal interactions between fronto-parietal central executive and default-mode networks in humans. *Proc. Natl Acad. Sci. USA* **110**, 19944–19949 (2013).
 8. Bandettini, P. A. Twenty years of functional MRI: the science and the stories. *Neuroimage* **62**, 575–588 (2012).
 9. Logothetis, N. K., Pauls, J., Augath, M., Trinath, T. & Oeltermann, A. Neurophysiological investigation of the basis of the fMRI signal. *Nature* **412**, 150–157 (2001).
 10. Attwell, D. *et al.* Glial and neuronal control of brain blood flow. *Nature* **468**, 232–243 (2010).
 11. Hillman, E. M. C. Coupling mechanism and significance of the bold signal: a status report. *Annu. Rev. Neurosci.* **37**, 161–181 (2014).
 12. Sirotni, Y. B. & Das, A. Anticipatory haemodynamic signals in sensory cortex not predicted by local neuronal activity. *Nature* **457**, 475–479 (2009).
 13. Thomsen, K., Offenhauser, N. & Lauritzen, M. Principal neuron spiking: neither necessary nor sufficient for cerebral blood flow in rat cerebellum. *J. Physiol. (Lond.)* **560**, 181–189 (2004).
 14. Einevoll, G. T., Kayser, C., Logothetis, N. K. & Panzeri, S. Modelling and analysis of local field potentials for studying the function of cortical circuits. *Nature Rev. Neurosci.* **14**, 770–785 (2013).
 15. Farah, M. J. Brain images, babies, and bathwater: critiquing critiques of functional neuroimaging. *Hastings Cent. Rep.* **44**, S19–S30 (2014).
 16. Poldrack, R. A. & Yarkoni, T. From brain maps to cognitive ontologies: informatics and the search for mental structure. *Annu. Rev. Psychol.* <http://dx.doi.org/10.1146/annurev-psych-122414-033729> (2015).
 17. Norman, K. A., Polyn, S. M., Detre, G. J. & Haxby, J. V. Beyond mind-reading: multi-voxel pattern analysis of fMRI data. *Trends Cogn. Sci.* **10**, 424–430 (2006).
 18. Kriegeskorte, N., Mur, M. & Bandettini, P. Representational similarity analysis—connecting the branches of systems neuroscience. *Frontiers Syst Neurosci* **2**, 4 (2008).
 19. Naselaris, T., Prenger, R. J., Kay, K. N., Oliver, M. & Gallant, J. L. Bayesian reconstruction of natural images from human brain activity. *Neuron* **63**, 902–915 (2009).
 20. Cowen, A. S., Chun, M. M. & Kuhl, B. A. Neural portraits of perception: reconstructing face images from evoked brain activity. *Neuroimage* **94**, 12–22 (2014).
 21. Mitchell, T. M. *et al.* Predicting human brain activity associated with the meanings of nouns. *Science* **320**, 1191–1195 (2008).
- An outstanding example of the use of fMRI along with a model of word meaning derived from a large text corpus to predict activation patterns associated with words.**
22. Huth, A. G., Nishimoto, S., Vu, A. T. & Gallant, J. L. A continuous semantic space describes the representation of thousands of object and action categories across the human brain. *Neuron* **76**, 1210–1224 (2012).
 23. Sprague, T. C., Saproo, S. & Serences, J. T. Visual attention mitigates information loss in small- and large-scale neural codes. *Trends Cogn. Sci.* **19**, 215–226 (2015).
 24. Lewis-Peacock, J. A. & Norman, K. A. Competition between items in working memory leads to forgetting. *Nature Commun.* **5**, 5768 (2014).
 25. Charest, I., Kievit, R. A., Schmitz, T. W., Deca, D. & Kriegeskorte, N. Unique semantic space in the brain of each beholder predicts perceived similarity. *Proc. Natl Acad. Sci. USA* **111**, 14565–14570 (2014).
 26. Davis, T. & Poldrack, R. A. Quantifying the internal structure of categories using a neural typicality measure. *Cereb. Cortex* **24**, 1720–1737 (2014).
 27. Mack, M. L., Preston, A. R. & Love, B. C. Decoding the brain's algorithm for categorization from its neural implementation. *Curr. Biol.* **23**, 2023–2027 (2013).
 28. Xue, G. *et al.* Greater neural pattern similarity across repetitions is associated with better memory. *Science* **330**, 97–101 (2010).
 29. Kriegeskorte, N. *et al.* Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron* **60**, 1126–1141 (2008).
- This paper applies representational similarity analysis to human fMRI and monkey electrophysiology data to demonstrate similar representational spaces in the inferior temporal cortex across species.**
30. Wager, T. D. *et al.* An fMRI-based neurologic signature of physical pain. *N. Engl. J. Med.* **368**, 1388–1397 (2013).
 31. Davis, T. *et al.* What do differences between multi-voxel and univariate analysis mean? How subject-, voxel-, and trial-level variance impact fMRI analysis. *Neuroimage* **97**, 271–283 (2014).
 32. Todd, M. T., Nystrom, L. E. & Cohen, J. D. Confounds in multivariate pattern analysis: theory and rule representation case study. *Neuroimage* **77**, 157–165 (2013).
 33. Dubois, J., de Berker, A. O. & Tsao, D. Y. Single-unit recordings in the macaque face patch system reveal limitations of fMRI MVPA. *J. Neurosci.* **35**, 2791–2802 (2015).
 34. Schultz, W. Predictive reward signal of dopamine neurons. *J. Neurophysiol.* **80**, 1–27 (1998).
 35. McClure, S. M., Berns, G. S. & Montague, P. R. Temporal prediction errors in a passive learning task activate human striatum. *Neuron* **38**, 339–346 (2003).
 36. O'Doherty, J. P., Dayan, P., Friston, K., Critchley, H. & Dolan, R. J. Temporal difference models and reward-related learning in the human brain. *Neuron* **38**, 329–337 (2003).
 37. Badre, D. & Frank, M. J. Mechanisms of hierarchical reinforcement learning in cortico-striatal circuits 2: evidence from fMRI. *Cereb. Cortex* **22**, 527–536 (2012).

38. Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P. & Dolan, R. J. Model-based influences on humans' choices and striatal prediction errors. *Neuron* **69**, 1204–1215 (2011).
39. LaRocque, K. F. *et al.* Global similarity and pattern separation in the human medial temporal lobe predict subsequent memory. *J. Neurosci.* **33**, 5466–5474 (2013).
40. Raichle, M. E. The brain's default mode network. *Annu. Rev. Neurosci.* **38**, 433–447 (2015).
41. Shulman, G. L. *et al.* Common blood flow changes across visual tasks: II. decreases in cerebral cortex. *J. Cogn. Neurosci.* **9**, 648–663 (1997).
42. Greicius, M. D., Krasnow, B., Reiss, A. L. & Menon, V. Functional connectivity in the resting brain: a network analysis of the default mode hypothesis. *Proc. Natl Acad. Sci. USA* **100**, 253–258 (2003).
43. Vincent, J. L. *et al.* Intrinsic functional architecture in the anaesthetized monkey brain. *Nature* **447**, 83–86 (2007).
44. Becerra, L., Pendse, G., Chang, P.-C., Bishop, J. & Borsook, D. Robust reproducible resting state networks in the awake rodent brain. *PLoS ONE* **6**, e25701 (2011).
45. Buckner, R. L. *et al.* Molecular, structural, and functional characterization of Alzheimer's disease: evidence for a relationship between default activity, amyloid, and memory. *J. Neurosci.* **25**, 7709–7717 (2005).
- This paper presents a multimodal analysis implicating the default mode network in cognitive decline associated with Alzheimer disease.**
46. Smith, S. M. *et al.* Correspondence of the brain's functional architecture during activation and rest. *Proc. Natl Acad. Sci. USA* **106**, 13040–13045 (2009).
- This paper demonstrates that resting state networks are systematically associated with cognitive functions.**
47. Laumann, T. O. *et al.* Functional system and areal organization of a highly-sampled individual human brain. *Neuron* **87**, 657–670 (2015).
48. Murphy, K., Birn, R. M. & Bandettini, P. A. Resting-state fMRI confounds and cleanup. *Neuroimage* **80**, 349–359 (2013).
49. Power, J. D., Schlaggar, B. L. & Petersen, S. E. Recent progress and outstanding issues in motion correction in resting state fMRI. *Neuroimage* **105**, 536–551 (2015).
50. Tagliazucchi, E. & Laufs, H. Decoding wakefulness levels from typical fMRI resting-state data reveals reliable drifts between wakefulness and sleep. *Neuron* **82**, 695–708 (2014).
51. Morcom, A. M. & Fletcher, P. C. Does the brain have a baseline? Why we should be resisting a rest. *Neuroimage* **37**, 1073–1082 (2007).
52. Nestler, E. J. & Hyman, S. E. Animal models of neuropsychiatric disorders. *Nature Neurosci.* **13**, 1161–1169 (2010).
53. Insel, T. R. The NIMH research domain criteria (RDoC) project: precision medicine for psychiatry. *Am. J. Psychiatry* **171**, 395–397 (2014).
54. Mayberg, H. S. Targeted electrode-based modulation of neural circuits for depression. *J. Clin. Invest.* **119**, 717–725 (2009).
55. Sulzer, J. *et al.* Real-time fMRI neurofeedback: progress and challenges. *Neuroimage* **76**, 386–399 (2013).
56. Goodkind, M. *et al.* Identification of a common neurobiological substrate for mental illness. *JAMA Psychiatry* **72**, 305–315 (2015).
- This paper examines a large structural imaging dataset and finds that brain abnormalities linked to mental illness are shared across diagnostic categories.**
57. Schizophrenia Working Group of the Psychiatric Genomics Consortium. Biological insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421–427 (2014).
58. Sperling, R. A., Karlawish, J. & Johnson, K. A. Preclinical Alzheimer disease—the challenges ahead. *Nature Rev. Neurol.* **9**, 54–58 (2013).
59. Owen, A. M. Detecting consciousness: a unique role for neuroimaging. *Annu. Rev. Psychol.* **64**, 109–133 (2013).
60. Ariely, D. & Berns, G. S. Neuromarketing: the hope and hype of neuroimaging in business. *Nature Rev. Neurosci.* **11**, 284–292 (2010).
61. Berns, G. S. & Moore, S. A neural predictor of cultural popularity. *J. Consum. Psychol.* **22**, 154–160 (2012).
62. Falk, E. B., Berkman, E. T. & Lieberman, M. D. From neural responses to population behavior: neural focus group predicts population-level media effects. *Psychol. Sci.* **23**, 439–445 (2012).
63. Gabrieli, J. D. E., Ghosh, S. S. & Whitfield-Gabrieli, S. Prediction as a humanitarian and pragmatic contribution from human cognitive neuroscience. *Neuron* **85**, 11–26 (2015).
64. Morse, S. J. Brain overclaim syndrome and criminal responsibility: A diagnostic note. *Ohio State J. Criminal Law* **3**, 397–412 (2006).
65. Jones, O. D., Wagner, A. D., Faigman, D. L. & Raichle, M. E. Neuroscientists in court. *Nature Rev. Neurosci.* **14**, 730–736 (2013).
66. Farah, M. J., Hutchinson, J. B., Phelps, E. A. & Wagner, A. D. Functional MRI-based lie detection: scientific and societal challenges. *Nature Rev. Neurosci.* **15**, 123–131 (2014).
67. Reardon, S. Neuroscience in court: the painful truth. *Nature* **518**, 474–476 (2015).
68. Underwood, E. DARPA aims to rebuild brains. *Science* **342**, 1029–1030 (2013).
69. Gerits, A. *et al.* Optogenetically induced behavioral and functional network changes in primates. *Curr. Biol.* **22**, 1722–1726 (2012).
70. Insel, T. R., Landis, S. C. & Collins, F. S. The NIH Brain Initiative. *Science* **340**, 687–688 (2013).
71. Moeller, S. *et al.* Multiband multislice GE-EPI at 7 tesla, with 16-fold acceleration using partial parallel imaging with application to high spatial and temporal whole-brain fMRI. *Magn. Reson. Med.* **63**, 1144–1153 (2010).
72. Yacoub, E., Harel, N. & Ugurbil, K. High-field fMRI unveils orientation columns in humans. *Proc. Natl Acad. Sci. USA* **105**, 10607–10612 (2008).
73. Magnain, C. *et al.* Optical coherence tomography visualizes neurons in human entorhinal cortex. *Neurophotonics* **2**, 015004 (2015).
- A demonstration of the power of optical coherence tomography to image neural structure in ex vivo human brain tissue.**
74. Hawrylycz, M. J. *et al.* An anatomically comprehensive atlas of the adult human brain transcriptome. *Nature* **489**, 391–399 (2012).
75. Richiardi, J. *et al.* Correlated gene expression supports synchronous activity in brain networks. *Science* **348**, 1241–1244 (2015).
76. Van Essen, D. C. *et al.* The WU-Minn Human Connectome Project: an overview. *Neuroimage* **80**, 62–79 (2013).
- This paper presents a broad overview of the Human Connectome Project.**
77. Saygin, Z. M. *et al.* Tracking the roots of reading ability: white matter volume and integrity correlate with phonological awareness in prereading and early-reading kindergarten children. *J. Neurosci.* **33**, 13251–13258 (2013).
78. Van Essen, D. C. *et al.* Mapping Connections in Humans and Non-human Primates: Aspirations and Challenges for Diffusion Imaging 2nd edn, Ch. 16 (Elsevier, 2013).
79. Eveley, C. *et al.* Superficial white matter fiber systems impede detection of long-range cortical connections in diffusion MR tractography. *Proc. Natl Acad. Sci. USA* **112**, E2820–E2828 (2015).
80. Smith, S. M. *et al.* Network modelling methods for fMRI. *Neuroimage* **54**, 875–891 (2011).
81. Honey, C. J. *et al.* Predicting human resting-state functional connectivity from structural connectivity. *Proc. Natl Acad. Sci. USA* **106**, 2035–2040 (2009).
82. Shen, K. *et al.* Information processing architecture of functionally defined clusters in the macaque cortex. *J. Neurosci.* **32**, 17465–17476 (2012).
83. Yarkoni, T., Poldrack, R. A., Nichols, T. E., Van Essen, D. C. & Wager, T. D. Large-scale automated synthesis of human functional neuroimaging data. *Nature Methods* **8**, 665–670 (2011).
84. Ioannidis, J. P. A. Why most published research findings are false: author's reply to Goodman and Greenland. *PLoS Med.* **4**, e215 (2007).
85. Button, K. S. *et al.* Power failure: why small sample size undermines the reliability of neuroscience. *Nature Rev. Neurosci.* **14**, 365–376 (2013).
86. Carp, J. On the plurality of (methodological) worlds: estimating the analytic flexibility of fMRI experiments. *Frontiers Neurosci.* **6**, 149 (2012).
87. Kriegeskorte, N., Simmons, W. K., Bellgowan, P. S. F. & Baker, C. I. Circular analysis in systems neuroscience: the dangers of double dipping. *Nature Neurosci.* **12**, 535–540 (2009).
88. Stein, J. L. *et al.* Identification of common variants associated with human hippocampal and intracranial volumes. *Nature Genet.* **44**, 552–561 (2012).
- This paper presents the first well-powered genome-wide association study of brain structure.**
89. Waskom, M. L., Kumaran, D., Gordon, A. M., Rissman, J. & Wagner, A. D. Frontoparietal representations of task context support the flexible control of goal-directed cognition. *J. Neurosci.* **34**, 10743–10755 (2014).
90. Irimia, A., Chambers, M. C., Torgerson, C. M. & Van Horn, J. D. Circular representation of human cortical networks for subject and population-level connectomic visualization. *Neuroimage* **60**, 1340–1351 (2012).
91. Blokland, G. A. M. *et al.* Heritability of working memory brain activation. *J. Neurosci.* **31**, 10882–10890 (2011).
92. Glahn, D. C. *et al.* Genetic control over the resting brain. *Proc. Natl Acad. Sci. USA* **107**, 1223–1228 (2010).
93. Barnett, J. H., Scoriels, L. & Munafò, M. R. Meta-analysis of the cognitive effects of the catechol-O-methyltransferase gene Val158/108Met polymorphism. *Biol. Psychiatry* **64**, 137–144 (2008).
94. Nickl-Jockschat, T., Janouschek, H., Eickhoff, S. B. & Eickhoff, C. R. Lack of meta-analytic evidence for an impact of COMT Val158Met genotype on brain activation during working memory tasks. *Biol. Psychiatry* <http://dx.doi.org/10.1016/j.biopsych.2015.02.030> (2015).
95. Farah, M. J. Neuroethics: the ethical, legal, and societal impact of neuroscience. *Annu. Rev. Psychol.* **63**, 571–591 (2012).
96. Illes, J. & Racine, E. Imaging or imagining? A neuroethics challenge informed by genetics. *Am. J. Bioeth.* **5**, 5–18 (2005).
97. Farah, M. J. & Gillihan, S. J. The puzzle of neuroimaging and psychiatric diagnosis: technology and nosology in an evolving discipline. *AJOB Neurosci.* **3**, 31–41 (2012).
98. Conrad, P. *The Medicalization of Society: On the Transformation of Human Conditions into Treatable Disorders* (Johns Hopkins Univ. Press, 2007).
99. Sahakian, B. & Morein-Zamir, S. Professor's little helper. *Nature* **450**, 1157–1159 (2007).
100. Fitz, N. S. & Reiner, P. B. The challenge of crafting policy for do-it-yourself brain stimulation. *J. Med. Ethics* **41**, 410–412 (2015).
101. Horvath, J. C., Forte, J. D. & Carter, O. Quantitative review finds no evidence of cognitive effects in healthy populations from single-session transcranial direct current stimulation (tDCS). *Brain Stimul.* **8**, 535–550 (2015).

Acknowledgements Thanks to I. Eisenberg, D. Glahn, R. Raizada, and M. Shine for comments on an earlier draft of this manuscript, and N. Logothetis for helpful discussions.

Author Contributions R.P. and M.F. planned and wrote the paper.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to R.P. (poldrack@stanford.edu).

A Cretaceous eutriconodont and integument evolution in early mammals

Thomas Martin¹, Jesús Marugán-Lobón^{2,3}, Romain Vullo⁴, Hugo Martín-Abad², Zhe-Xi Luo⁵ & Angela D. Buscalioni²

The Mesozoic era (252–66 million years ago), known as the domain of dinosaurs, witnessed a remarkable ecomorphological diversity of early mammals. The key mammalian characteristics originated during this period and were prerequisite for their evolutionary success after extinction of the non-avian dinosaurs 66 million years ago. Many ecomorphotypes familiar to modern mammal fauna evolved independently early in mammalian evolutionary history. Here we report a 125-million-year-old eutriconodontan mammal from Spain with extraordinary preservation of skin and pelage that extends the record of key mammalian integumentary features into the Mesozoic era. The new mammalian specimen exhibits such typical mammalian features as pelage, mane, pinna, and a variety of skin structures: keratinous dermal scutes, protospines composed of hair-like tubules, and compound follicles with primary and secondary hairs. The skin structures of this new Mesozoic mammal encompass the same combination of integumentary features as those evolved independently in other crown Mammalia, with similarly broad structural variations as in extant mammals. Soft tissues in the thorax and abdomen (alveolar lungs and liver) suggest the presence of a muscular diaphragm. The eutriconodont has molariform tooth replacement, ossified Meckel's cartilage of the middle ear, and specialized xenarthrous articulations of posterior dorsal vertebrae, convergent with extant xenarthran mammals, which strengthened the vertebral column for locomotion.

Mammals have diverse integumentary structures such as hair, spines, and scutes¹. Fossilized fur is known for Mesozoic mammals from the Middle Jurassic period onwards, and can be traced back to the early-divergent mammaliaform clades^{2–7}. So far, however, fur has fossilized only as hair impressions and compressions, types of preservation that rarely show microstructure. The exceptional phosphatized preservation of skin tissues in the specimen from Las Hoyas extends direct knowledge of mammalian integumentary microstructure from 60 million years ago (Ma) (Late Palaeocene epoch)⁸ back to 125–127 Ma (Early Cretaceous). The new Spanish gobiconodontid (Eutriconodonta, Mammalia) combines a complete and articulated skeleton, with extraordinary preservation of skin, hair, keratinous dermal scutes, and remnants of visceral organs.

The most striking feature of extant mammalian hairs is their outstanding polymorphism, among closely related species, at different developmental stages, and even in the same individual⁹. The new fossil shows the prominent polymorphism in a variety of keratinous appendages, some of which are described for Mesozoic mammals for the first time. The cylindrical structures (protospines) are formed from the fusion of cylindrical micro-tubules, architecturally similar to the spines of modern mammals. Such protospines represent the earliest record of a microstructure analogous to that of spiny mice (*Acomys*) and hedgehogs (*Erinaceus*). As in extant mammals, hair shafts have three architectural units: cortex, medulla, and cuticles with varied designs. Moreover, at a microscopic scale the epidermal layer preserves compound follicles with more than two hairs stemming from the same follicle. The structural diversity of hairs in this new mammaliaform indicates that polymorphism of hairs and related structures is probably correlated with the same developmental process as in extant mammals⁹.

On its posterior dorsal vertebrae the new eutriconodont has additional xenarthrous articulations in addition to the usual zygapophyses. This is well known for Xenarthra (armadillos, anteaters, and sloths) among

extant placentals, and similar structures have been observed in some Eocene eutherians^{10–12}. The fact that features of xenarthral articulation are also present in the Late Jurassic mammaliaform *Fruitafossor*¹³ makes the new mammal from Spain another case of convergent evolution of this remarkable skeletal feature that strengthens the dorsal vertebral column for versatile locomotor functions.

Class Mammalia¹⁴

Order Eutriconodonta¹⁵

Family Gobiconodontidae¹⁶

Spinolestes xenarthrosus gen. et sp. nov.

Etymology. *Spinosus* (Latin), in reference to the spiny integument; *λέστης* (Greek) or *lestes* (Latin spelling), meaning robber and a common term in taxonomic names of mammals. The specific name *xenarthrosus* refers to the special additional (ξένος, (Greek) strange) articulation facets (ἄρθρον, (Greek) articulation) of the dorsal vertebrae.

Holotype. MCCMLH30000, Museo de las Ciencias de Castilla-La Mancha. A complete skeleton with integumentary structures preserved. Skeletal and integumentary remains of the slab MCCMLH30000A were transferred to a matrix of epoxy resin in preparation, to expose the embedded side of the fossil, but are kept intact in the counter slab MCCMLH30000B (Figs 1 and 2a and Extended Data Fig. 1). Life Science Identifier: urn:lsid:zoobank.org:act:10B31072-A727-4663-A7DC-4224FB97E1E3.

Locality and horizon. Las Hoyas Quarry, Calizas de la Huérgina Formation, southwestern Iberian Basin (Cuenca, Spain). Las Hoyas is latest Barremian (125–127 Ma) in age, on the basis of charophytes and ostracodes¹⁷. The Las Hoyas Konservat-Lagerstätte occurs in finely laminated limestones deriving from a freshwater wetland. Fossils are usually preserved fully articulated, including soft tissues such as mineralized muscle and skin. Potential mechanisms for

¹Steinmann-Institut für Geologie, Mineralogie und Paläontologie, Universität Bonn, Nussallee 8, 53115 Bonn, Germany. ²Unidad de Paleontología, Departamento de Biología, Facultad Ciencias, Universidad Autónoma de Madrid, C/Darwin 2, 28049 Madrid, Spain. ³Dinosaur Institute, Natural History Museum of Los Angeles County, 900 Exposition Boulevard, Los Angeles, California 90007, USA. ⁴Géosciences Rennes, UMR CNRS 6118, Université de Rennes 1, Campus de Beaulieu, bâtiment 15, 263 avenue du Général Leclerc, CS 74205, 35042 Rennes, France. ⁵Department of Organismal Biology and Anatomy, The University of Chicago, 1027 East 57th Street, Chicago, Illinois 60637, USA.

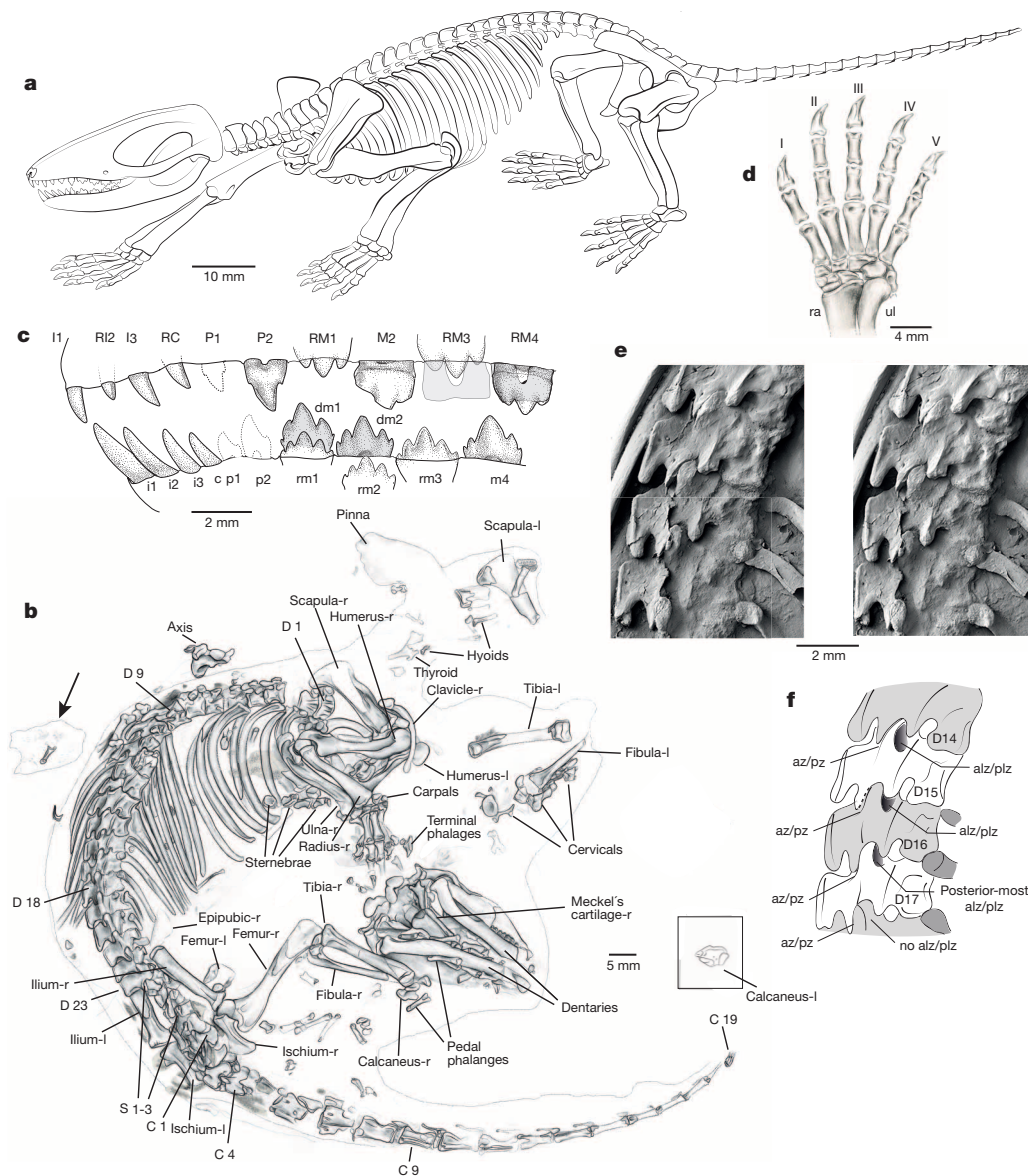


Figure 1 | New Early Cretaceous gobiconodontid *Spinolestes xenarthrosus*.

a, Skeletal reconstruction from the type specimen. **b**, Drawing of holotype MCCMLH30000A with skeletal element identification and outline of preserved integument. Arrow points to skin patch preserving hair bundles, placed between dorsal vertebrae 11 and 14. Inset: left calcaneus in dorsal view, on the transferred slab separate from the skeleton. Abbreviations: C, caudal vertebra; D, dorsal vertebra; l, left; r, right; S, sacral vertebra; c, Reconstruction of dentition with deciduous molariforms (grey shading). Abbreviations: C/c,

upper/lower canine; I/i, incisor; M/m, molariform; P/p, premolar; D/d, deciduous tooth; R/r, replacing tooth. **d**, Reconstruction of right hand. Numerals designate finger rays. Abbreviations: ra, radius; ul, ulna. **e**, Dorsal vertebrae 14–17 with xenarthrous articulations, SEM images (stereo pair). **f**, Drawing of dorsal vertebrae 14–17 with accessory (xenarthrous) lateral articulations. Abbreviations: az, anterior zygapophysis; alz, anterior lateral zygapophysis; pz, posterior zygapophysis; plz, posterior lateral zygapophysis.

exquisite preservation are microbial mats, anoxia, and rapid burial by sediments¹⁸.

Diagnosis. Dental formula $I^3-C^1-P^2-M^4/i_3-c_1-p_2-m_4$. Cheek teeth of general ‘amphilestid’ pattern (Fig. 1c and Extended Data Fig. 2a–c). Differs from stem mammaliaforms in lacking the postdentary trough. Differs from the ‘amphilestids’ *Amphilestes* and *Phascolotherium* in fewer premolars, fewer molariforms, or both; from *Hakusanodon*¹⁹ and *Juchilestes*²⁰ in fewer molariforms. Differs from *Jeholodens* and *Yanoconodon*^{21,22} by the presence of accessory cusps d and e of lower molariforms and a lingual cingulum in upper molariforms. Differs from *Liaconodon*²³ in having one fewer lower incisor and one more upper molariform, and lacking the uniform enlargement of incisors and canines of the latter. Differs from genera of the Triconodontidae in absence of vertical tongue-in-groove interlock of molariforms, and more uniform heights of cusps of the latter. Unique among gobiconodontids

in that the erupting (replacing) molariform is lingual (side-by-side) to the deciduous molariform of the same tooth locus in the maxilla (Supplementary Video 1). Within gobiconodontids, *Spinolestes* differs from *Gobiconodon*^{24–26} in having more incisors and fewer premolars. Of all gobiconodontids, *Spinolestes* is most similar to *Repenomamus* in dentition, but differs from *Repenomamus* species in having one fewer molariform each for uppers and for lowers, in enlargement of lower first incisor, and in having both labial and lingual cingula in upper molariforms²⁷ (for full diagnosis and hypotheses of phylogenetic placement see Extended Data Fig. 3 and Supplementary Information).

Description

The skeleton is exposed on its right lateral side in the main slab after its transfer to resin embedding. The skull is exposed on its ventral side (Figs 1b and 2a and Extended Data Fig. 1) and has a rounded rostrum.

The jugals are complete, extending from anterior zygoma to the squamosal glenoid. The glenoid is transversely expanded. The petrosal has a cylindrical promontorium, and a narrow lateral trough. The occipital condyles are large and oval shaped. The mandibles are robust for their short length, and have no angular process. The dentary condyles are broadly oval and bent laterally from the plane of coronoid process and mandibular body. Meckel's cartilage is ossified, preserved *in situ* on the left mandible but detached on the right side, exposing a wide Meckel's groove on the right mandible (Supplementary Video 2). The middle ear bones themselves are lost, but given the curved Meckel's cartilage, we infer the ear was medio-laterally separate from the mandible yet still connected antero-posteriorly to the mandible via the ossified Meckel's cartilage, as in *Yanoconodon*, *Liaoconodon*, and other gobiconodontids^{22,23,25,28}.

Spinolestes shows lower molariform replacements at least at m1 to m2 loci (Fig. 1c and Extended Data Fig. 2a). The right m1 and m3 are in the process of eruption, whereas beneath deciduous m2 the germ of replacing m2 is visible. Upper molariform replacements occur from M1 to M4 positions. The right M1 is in the process of eruption. Deciduous M3 and M4 are heavily worn in an oblique angle, down to the roots, but still in place in the maxilla, while replacing M3 and M4 are erupting to the lingual side of DM3, DM4 although at different stages (Supplementary Video 1). Medial molariform replacement/succession occurs only in the maxilla, but mandible shows vertical replacement of deciduous by erupting molariforms. Because upper M2 appears to have erupted earlier than M1 and M3, and lower m1 and m3 erupted before m2, *Spinolestes* shows alternative replacements for these tooth loci. Alternating replacement is common in basal mammals²⁹, and is consistent with replacement patterns of anterior molariforms in *Gobiconodon ostromi*²⁴.

The shoulder girdle is therian-like (Fig. 1a, b). The broad, triangular scapula has a high spine that ends in a cranially oriented acromion. Its posterior border is laterally bent, forming a secondary spine. As in therians, the glenoid fossa is oval, oriented nearly perpendicular to the scapular plate. The curved clavicle has a large, curved synovial facet on its distal end for mobile articulation with the acromion. The coracoid is small and fused to the scapula. There is no procoracoid. Five ossified sternal elements are present, plus the partially preserved interclavicle in the ventral part of the thorax. The claviculo-interclavicular joint is mobile, as documented for other eutriconodonts^{21,22,27}.

Spinolestes has 16 thoracic vertebrae, seven lumbar vertebrae (defined by zygapophyses rotated to vertical, and by transverse processes), three fused sacra, and 22 caudals. Twenty dorsal vertebrae bear ribs, of which 16 pairs are associated with the thoracic and four pairs with the lumbar vertebrae. Ultimate and penultimate lumbar have fused ribs. The robust ribs are two-headed, and have rounded cross sections. The dorsal vertebrae have reclined spinal processes. Beginning with dorsal vertebra 11, the spinal processes become larger in the antero-posterior dimension, and develop an expanded flat top. This flat top becomes larger and more prominent in the successively posterior lumbar.

Spinolestes shows accessory lateral xenarthrous articulations (*sensu* ref. 30) from dorsal vertebrae 9/10 to dorsals 16/17 (Fig. 1e, f and Extended Data Fig. 2d). The anterolateral aspect of the prezygapophysis forms a lateral zygapophyseal articulation with the preceding vertebra. Among therian mammal groups, xenarthrous articulation is most prominently developed in South American xenarthrans (armadillos, anteaters, and sloths), and in some Eocene placentals^{10–12}. Among Mesozoic mammaliaforms, a xenarthrous articulation is present in *Fruitafossor*¹³. Besides *Spinolestes*, this derived structure is also present in other gobiconodontids according to our observation (also

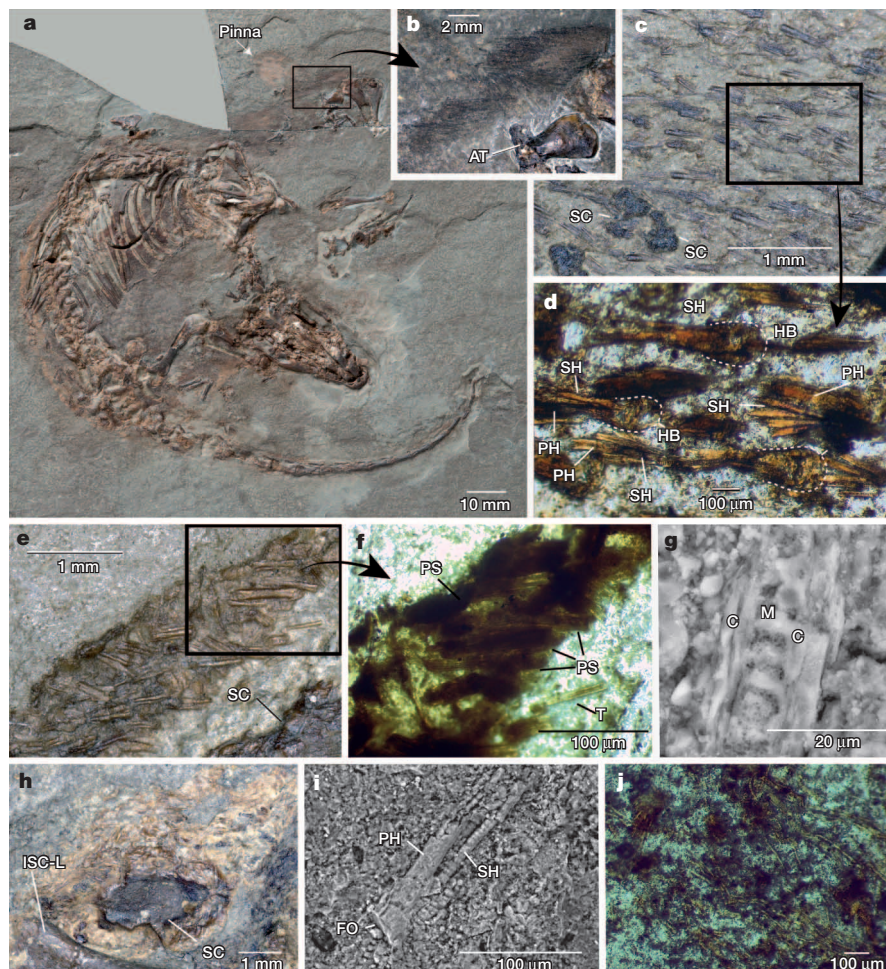


Figure 2 | New gobiconodontid *S. xenarthrosus*, holotype counter slab (MCCMLH30000B) and integumentary structures. **a**, Holotype counter slab with its half skeleton on original rock matrix. **b**, Guard hairs of dense fur (mane) in the cervical region (AT, atlas). **c**, Hairs and keratinous dermal scutes (SC) from portion of integument placed between dorsal vertebrae 11 and 14. **d**, Arrangement of compound hair follicles (inset) including hair bulbs (HB), primary hairs (PH), and secondary hairs (SH). **e**, Area with protospines, dermal scute (SC), and isolated tubules (inset). **f**, Protospines (PS) and isolated tubules (T) under translucent light. **g**, SEM image of a primary hair shaft's longitudinal section displaying inner structure with cortex (C) and medulla of discontinuous type (M). **h**, Oval keratinous dermal scute (SC) dorsally of left ischium (ISC-L) with tubules and the homogeneous matrix. **i**, SEM image of compound follicle (FO) with primary (thicker) hair of mosaic cuticular pattern and secondary (thinner) hairs with coronary (annulus) cuticular pattern. **j**, Underfur hairs with annular cuticular pattern in the abdominal region under translucent light. **c–f, j**, Taken from transferred slab MCCMLH30000A.

ref. 27: figures 3-8 and 3-10). However, gobiconodontids including *Spinolestes*, and the mammaliaform *Fruitafossor*, lack the fully developed mammillary process (metapophysis³¹) of the xenarthrous articulation seen in xenarthran placentals.

Phylogenetic relationships

Two phylogenetic analyses were conducted (Extended Data Fig. 3). One estimates the relationships of *Spinolestes* among major mammaliaform clades (Extended Data Fig. 3b on a data set of 111 taxa and 490 characters^{7,28}), and the other evaluates the placement of *Spinolestes* among 20 genera of eutriconodonts (Extended Data Fig. 3a, on the basis of 90 dental and mandibular characters^{19,20,32,33}). *Spinolestes* is nested within Eutriconodontia as a member of Gobiconodontidae. Within Gobiconodontidae, *Spinolestes* is the sister-taxon of the clade *Gobiconodon* + *Repenomamus*. The finer sampling of eutriconodont taxa improves the resolution of the known eutriconodont phylogenies^{19,32,33}. In fact, gobiconodontids have a broad distribution from Asia to North America during the Early Cretaceous (Barremian to Albian), and are relatively abundant in Asia. The discovery of *Spinolestes*, plus two gobiconodontid teeth from the Early Cretaceous of Spain³⁴ and Britain³⁵, has extended the distribution of this family to Western Europe. This is consistent with the broad interchanges of other mammal groups on the Laurasian landmass during the Early Cretaceous²⁹.

Integumentary and soft-tissue structures

The *Spinolestes* specimen shows regional hair differentiation and exquisite preservation of many integumentary structures (Fig. 2 and Extended Data Figs 1 and 4-8). A dense mane of long guard hairs (3-5 mm long) (Fig. 2b) is developed in the parietal, cervical, and scapular regions. Long and fine hairs are also present along the dorsal region, where they form a median crest, and along most of the tail. The rest of the body is covered by a soft and dense underfur (Fig. 2j).

Multiple pieces of skin are preserved with microstructural details of hair such as primary and gradually smaller secondary hairs that are organized in compound follicles (Fig. 2c, d, i and Extended Data Figs 4-6), as in many modern placentals³⁶⁻³⁸, marsupials, and monotremes³⁹. In an isolated piece of skin between dorsal vertebrae 11 and 14, compound follicles are associated with scale-like folds of the skin (Extended Data Fig. 4c), identical to the pattern of hairy skin of extant dogs³⁶ (Extended Data Fig. 4d). The cuticular scales have a smooth surface. The cuticula of primary hairs shows an irregular mosaic of imbricate, ovate scales (Fig. 2i and Extended Data Fig. 5a, b). The cuticula of secondary hairs displays coronal scales (annulus and encircling the shaft) with either simple or serrate free margins (Fig. 2i and Extended Data Fig. 5b, c). The longitudinal section of a primary hair shaft exhibits a discontinuous medulla with serial chambers and septa, and a thick cortex made of typical fusiform, spindle-shaped cells (ref. 40: figure 3) (Fig. 2g and Extended Data Fig. 6e). In a few regions of the body, the hairs have short, truncated shafts less than 1 mm long (Fig. 2c, d) with dark-coloured distal ends. Compared with normal hairs, such features match the so-called 'block hairs' and 'i-hairs' that represent a pathological condition associated with a fungal infection (dermatophytosis) widely spread in extant mammals (ref. 41: figure 2.50).

Spinolestes has what we call here protospines in the dorso-sacral region of the body (Fig. 2e, f), varying in basal diameters from 80 to 130 μm . Each of these small protospines is made up of two or more longitudinally fused tubules that appear to have scaly external surface, and an interior medullary cavity surrounded by a cortical structure (Fig. 2f). These tubules are interpreted as modified primary and secondary hairs. Most of the thicker protospines cluster with thinner ones and with isolated tubules, all of which are oriented randomly. A dozen or so oval dermal scutes up to 4 mm in dimension are present in the dorsal, lumbar, and pelvic regions (Fig. 2h and Extended Data Fig. 8). These scutes consist of a random arrangement of tubules merging to a homogeneous matrix.

Spines and hairs have fundamentally similar structural designs^{9,42}. In embryogenesis, spines appear later than hairs and are developed from the fusion of several hair follicles either by the merger of multiple hair follicles of similar size, or by one large merging with several smaller surrounding hair follicles⁹. Scanning electron microscopy (SEM) revealed that the protospines of *Spinolestes* are formed by the merger of several hair-like tubules, which is similar to the development of spines from the fusion of hair follicles in extant mammals. This suggests that the evolution of spines was mediated in Mesozoic mammals by a similar process of fusion of multiple hair follicles, as in the embryogenesis of integumentary structures in extant mammals. The unambiguous presence of compound hair follicles and their related scale-like skin folds in such a basal mammal clade as eutriconodonts confirm that these features are ancestral to mammals. *Spinolestes* further proves that mammalian underfur, guard hair, and spines had already fully differentiated in some Mesozoic mammals during the Early Cretaceous. The fact that multiple specimens of many taxa of eutriconodonts are preserved with fur but lack spines makes the presence of protospines in *Spinolestes* unique among eutriconodonts, apparently evolving independently from monotremes (echidnas), and some placentals such as hedgehogs and spiny mice. Although spines of the spiny mouse *Acomys* consist of a single modified awl hair, spiny mice can serve as a modern analogue to *Spinolestes* with spiny hairs located on the lower dorsum, possibly for display and protection by easy loss of spines when bitten by predators on the back⁴³.

The specimen also shows a unique preservation of several organs. In the area of the scalp, the left outer ear (external pinna) is perfectly preserved (Fig. 2a and Extended Data Fig. 7). Within the thoracic ribcage of *Spinolestes*, a patch of fossilized soft tissues contains tubular structures with a branching pattern (Extended Data Fig. 9). From the position and distribution in the rib cage, this most probably represents fossilized lung tissue, and the branching structures probably represent the bronchioles of the lung⁴⁴. Posteriorly to the lung tissue, a large oval area of reddish-brown soft-tissue (Extended Data Fig. 9a) is interpreted as residues of the liver according to its anatomical position and colour. Liver tissue is rich in iron, and provides a reddish colour, as has been reported for the theropod *Scipionyx* from the late Early Cretaceous (Albian) of Italy⁴⁵. The boundary between the inferred lung and the liver tissues extends, obliquely, from the distal tip of the 3rd rib to the proximal end of the 15th rib (Extended Data Fig. 9a), corresponding to the muscular diaphragm in extant mammals³⁸. This extends from the mid-sternal area to the posterior-most thoracic vertebrae, between the lungs in the thorax, and the liver in the abdominal cavity. Its presence reconfirms that this complex respiratory apparatus, which is tightly correlated with locomotion⁴⁶, was already functional in Mesozoic mammals.

Skeletal and locomotor adaptations

The scaling of long bones and mandible to body mass^{47,48} estimates that the weight of *Spinolestes* ranged from 52 to 72 g (see biometrics in Supplementary Information), indicating that it was a medium-sized Mesozoic mammal, falling within the range of small extant didelphid marsupials (for example, the size of *Monodelphis brevicaudata* and *Marmosa murina*)⁴⁷. Its limb proportions and xenarthrous vertebrae configuration together encompass a singular ecomorphological combination. The fibula is not co-ossified with the tibia, and is only slightly thinner than the latter. The carpals are relatively short, and the value of the phalangeal index⁴⁹ of digit III (equal to 121) corresponding to the top quartile of mammals with a terrestrial locomotory mode definitively indicates that it was not arboreal. The terminal phalanges lack a highly curved dorsal outline (Fig. 1d), but are robust and relatively wide, both of which are features common for extant mammals of terrestrial and fossorial lifestyles. Many extant placentals with xenarthrous vertebral articulations (or with reinforced thoracolumbar vertebrae in general) are capable diggers, such as armadillos and anteaters. The Jurassic mammaliaform *Fruitafossor* exhibits striking convergent adaptations to armadillos such as reduced enamel-less peglike teeth and xenarthrous

vertebral articulations. The Eocene mammal *Eurotamandua* (putatively a palaeonodont) has edentulous jaws and strongly enlarged manual claws. Both taxa are considered as fossorial and feeding on colonial insects^{10,13}. In contrast to *Fruitafossor*, *Eurotamandua*, and xenarthrans, the dentition of *Spinolestes* does not exhibit any tendency to reduction. Thus, *Spinolestes* shows a mosaic of functional features, suggesting a terrestrial locomotion, with ambulatory gait, as in *Gobiconodon*^{24,27}, and potential digging abilities, with a more versatile lifestyle (Extended Data Fig. 10). The extant armoured shrew *Scutisorex* can serve as a modern analogue for the lifestyle of *Spinolestes*. *Scutisorex* is unique among mammals by its massive interlocking lumbar vertebrae, which provide an extraordinary vertebral strength. It has been hypothesized that *Scutisorex* uses this vertebral strength to force open the base of palm leaves to search for insects or larvae in swampy forests⁵⁰. A similar lifestyle, supported by the strong arms and hands, is hypothesized here for *Spinolestes*, which lived in the vegetated Las Hoyas wetland environment.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 5 March; accepted 13 July 2015.

- Alibardi, L. Perspectives on hair evolution based on some comparative studies on vertebrate cornification. *J. Exp. Zool.* **318B**, 325–343 (2012).
- Ji, Q. *et al.* The earliest known eutherian mammal. *Nature* **416**, 816–822 (2002).
- Ji, Q. *et al.* A swimming mammaliaform from the Middle Jurassic and ecomorphological diversification of early mammals. *Science* **311**, 1123–1127 (2006).
- Meng, Q.-J. *et al.* An arboreal docodont from the Jurassic and mammaliaform ecological diversification. *Science* **347**, 764–768 (2015).
- Meng, J. *et al.* A Mesozoic gliding mammal from northeastern China. *Nature* **444**, 889–893 (2006).
- Vullo, R. *et al.* Mammalian hairs in Early Cretaceous amber. *Naturwissenschaften* **97**, 683–687 (2010).
- Zhou, C.-F. *et al.* A Jurassic mammaliaform and the earliest mammalian evolutionary adaptations. *Nature* **500**, 163–167 (2013).
- Meng, J. & Wyss, A. R. Multituberculate and other mammal hair recovered from Palaeogene excreta. *Nature* **385**, 712–714 (1997).
- Chernova, O. F. Evolutionary aspects of hair polymorphism. *Biol. Bull.* **33**, 43–52 (2006).
- Storch, G. *Eurotamandua joresi*, ein Myrmecophagide aus dem Eozän der “Grube Messel” bei Darmstadt (Mammalia, Xenarthra). *Senckenbergiana Lethaea* **61**, 247–289 (1981).
- Ting, S.-Y. A Paleocene edentate from Nanxiong Basin, Guangdong. *Palaeontol. Sin. New Series C* **17**, 85–118 (1987).
- Rose, K. D. & Emry, R. J. in *Mammal Phylogeny: Placentals* (eds Szalay, F. S., Novacek, M. J. & McKenna, M. C.) 81–102 (Springer, 1993).
- Luo, Z.-X. & Wible, J. R. A Late Jurassic digging mammal and early mammalian diversification. *Science* **308**, 103–107 (2005).
- Linnaeus, C. *Systema Naturae per Regna Tria Naturae, Secundum Classes, Ordines, Genera, Species, cum Characteribus, Differentiis, Synonymis, Locis* Vol. 1 (Regnum Animale) editio decima, reformata (Laurentius Salvius, 1758).
- Kermack, K. A. *et al.* The lower jaw of *Morganucodon*. *Zool. J. Linn. Soc.* **53**, 87–175 (1973).
- Chow, M. & Rich, T. H. V. A new triconodontan (Mammalia) from the Jurassic of China. *J. Vertebr. Paleontol.* **3**, 226–231 (1984).
- Fregenal-Martínez, M. A. & Meléndez, N. in *Lake Basins through Space and Time* (eds Gierlowski-Kordesch, E. H. & Kelts, K.) AAPG Studies in Geology Vol. 46, 303–314 (2000).
- Buscalioni, A. D. & Fregenal-Martínez, M. A. A holistic approach to the palaeoecology of Las Hoyas Konservat-Lagerstätte (La Huérgina Formation, Lower Cretaceous, Iberian Ranges, Spain). *J. Iber. Geol.* **36**, 297–326 (2010).
- Rougier, G. W. *et al.* An Early Cretaceous mammal from the Kuwajima Formation (Tetori Group), Japan, and a reassessment of triconodont phylogeny. *Ann. Carnegie Mus.* **76**, 73–115 (2007).
- Gao, C.-L. *et al.* A new mammal skull from the Lower Cretaceous of China with implications for the evolution of obtuse-angled molars and ‘amphilestid’ eutriconodonts. *Proc. R. Soc. B* **277**, 237–246 (2010).
- Ji, Q. *et al.* A Chinese triconodont mammal and mosaic evolution of the mammalian skeleton. *Nature* **398**, 326–330 (1999).
- Luo, Z.-X. *et al.* A new eutriconodont mammal and evolutionary development in early mammals. *Nature* **446**, 288–293 (2007).
- Meng, J. *et al.* Transitional mammalian middle ear from a new Cretaceous Jehol eutriconodont. *Nature* **472**, 181–185 (2011).
- Jenkins, F. A. & Schaff, C. R. The Early Cretaceous mammal *Gobiconodon* (Mammalia, Triconodonta) from the Cloverly Formation in Montana. *J. Vertebr. Paleontol.* **8**, 1–24 (1988).
- Kielan-Jaworowska, Z. & Dashzeveg, D. Early Cretaceous amphilestid (“triconodont”) mammal from Mongolia. *Acta Palaeontol. Pol.* **43**, 413–438 (1998).
- Li, C.-K. *et al.* A new species of *Gobiconodon* (Triconodonta, Mammalia) and its implication for the age of Jehol Biota. *Chin. Sci. Bull.* **48**, 1129–1134 (2003).
- Hu, Y.-M. *Postcranial Morphology of Repenomamus (Eutriconodonta, Mammalia): Implications for the Higher-Level Phylogeny of Mammals*. Dissertation, City Univ. of New York (2006).
- Luo, Z.-X. Developmental patterns in Mesozoic evolution of mammal ears. *Annu. Rev. Ecol. Syst.* **42**, 355–380 (2011).
- Kielan-Jaworowska, Z. *et al.* *Mammals from the Age of Dinosaurs. Origins, Evolution, and Structure* (Columbia Univ. Press, 2004).
- Gaudin, T. J. The morphology of xenarthrous vertebrae (Mammalia: Xenarthra). *Fieldiana Geol. New Series* **41**, 1–38 (1999).
- Lessertisseur, J. & Saban, R. in *Traité de Zoologie* Tome 16, Fascicule 1 (Mammifères, Téguments et Squelette) (ed. Grassé, P. P.) 709–1078 (Masson et Cie, 1967).
- Kusuhashi, N. *et al.* New triconodontids (Mammalia) from the Lower Cretaceous Shale and Fuxin formations, northeastern China. *Geobios* **42**, 765–781 (2009).
- Gaetano, L. C. & Rougier, G. W. New materials of *Argentoconodon fariarum* (Mammaliaformes, Triconodontidae) from the Jurassic of Argentina and its bearing on triconodont phylogeny. *J. Vertebr. Paleontol.* **31**, 829–843 (2011).
- Cuenca-Bescós, G. & Canudo, J. I. A new gobiconodontid mammal from the Early Cretaceous of Spain and its palaeogeographic implications. *Acta Palaeontol. Pol.* **48**, 575–582 (2003).
- Sweetman, S. C. A gobiconodontid (Mammalia, Eutriconodonta) from the Early Cretaceous (Barremian) Wessex Formation of the Isle of Wight, southern Britain. *Palaeontology* **49**, 889–897 (2006).
- Lovell, J. E. & Getty, R. The hair follicle, epidermis, dermis, and skin glands of the dog. *Am. J. Vet. Res.* **18**, 873–885 (1957).
- Whiteley, H. J. Giant compound hair follicles in the skin of the rabbit. *Nature* **181**, 850 (1958).
- Evans, H. E. & de Lahunta, A. *Miller's Anatomy of the Dog* (Elsevier Saunders, 2012).
- Spencer, B. & Sweet, G. The structure and development of the hairs of monotreme and marsupials. *Q. J. Microsc. Sci.* **36**, 549–588 (1899).
- Hausman, L. A. The cortical fusi of mammalian hair shafts. *Am. Nat.* **66**, 461–470 (1932).
- Rudnicka, L. *et al.* in *Atlas of Trichoscopy: Dermoscopy in Hair and Scalp Disease* (eds Rudnicka, L., Olszewska, M. & Rakowska, A.) 11–46 (Springer, 2012).
- Vincent, J. F. V. & Owers, P. Mechanical design of hedgehog spines and porcupine quills. *J. Zool.* **210**, 55–75 (1986).
- Montandon, S. A. *et al.* Two waves of anisotropic growth generate enlarged follicles in the spiny mouse. *EvoDevo* **5**, 33 (2014).
- Weibel, E. R. *et al.* Design of peripheral airways for efficient gas exchange. *Respir. Physiol. Neurobiol.* **148**, 3–21 (2005).
- Dal Sasso, C. & Signore, M. Exceptional soft tissue preservation in a theropod dinosaur from Italy. *Nature* **392**, 383–387 (1998).
- Bramble, D. M. & Jenkins, F. A. Mammalian locomotor-respiratory integration: implications for diaphragmatic and pulmonary design. *Science* **262**, 235–240 (1993).
- Campione, N. E. & Evans, D. C. A universal scaling relationship between body mass and proximal limb bone dimensions in quadrupedal terrestrial tetrapods. *BMC Biol.* **10**, 60 (2012).
- Foster, J. R. Preliminary body mass estimates for mammalian genera of the Morrison Formation (Upper Jurassic, North America). *PaleoBios* **28**, 114–122 (2009).
- Kirk, E. C. *et al.* Intrinsic hand proportions of eumarchontans and other mammals: Implications for the locomotor behavior of plesiadapiforms. *J. Hum. Evol.* **55**, 278–299 (2008).
- Stanley, W. T. *et al.* A new hero emerges: another exceptional mammalian spine and its potential adaptive significance. *Biol. Lett.* **9**, 20130486 (2013).

Supplementary Information is available in the online version of the paper.

Acknowledgements Research funds were provided by Spanish MINECO, Project CGL-2013-42643 P and Junta de Comunidades de Castilla-La Mancha. We thank J. L. Sañudo for finding the specimen, M. Llandres and O. Dülfer for preparation, D. Kranz for artwork, O. Sanisidro for the lifelike reconstruction of *S. xenarthrosus*, G. Oleschinski for photography, M. Furió and A. Valera for SEM, K. Jäger for three-dimensional reconstructions, and T. McCann for improving the English. R. L. Cifelli is thanked for review and B. Mähler for discussion on actuo-taphonomical experiments on rat carcasses.

Author Contributions A.D.B. designed the Las Hoyas research project; J.M.-L., R.V., H.M.-A., and A.D.B. participated in the fieldwork; T.M., J.M.-L., R.V., H.M.-A., Z.-X.L., and A.D.B. organized and conducted the research (preparation, computed tomography scan, light microscopy, and SEM imaging) and analysed data; Z.-X.L. performed phylogenetic analyses; and T.M., J.M.-L., R.V., Z.-X.L., and A.D.B. wrote the manuscript.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to T.M. (tmartin@uni-bonn.de), Z.-X.L. (zxluo@uchicago.edu) or A.D.B. (angela.delgado@uam.es).

METHODS

Upon discovery, the rock containing the specimen was split into two slabs, each containing original bone, fossilized soft tissues, and natural moulds of bones. One slab (MCCMLH30000A) was transferred to a matrix of epoxy resin to expose undisturbed bone surface hidden in the rock and to obtain casts of the natural moulds of bone impressions. After embedding in epoxy resin, the limestone matrix was removed by a formic acid bath that dissolved the calcium carbonate but not the hydroxyapatite of bone and the phosphatized soft tissue. The transparent artificial plastic matrix allows examination of integument structures under translucent light. The other slab (MCCMLH30000B) was left untouched in the limestone matrix to keep its original embedding situation.

The transferred slab and the non-transferred slab were scanned on a v|tome|x s μ CT-scanner at 240 kV (GE Electronics Phoenix|X-ray), 864 slices at a voxel size of 107.73 μ m. Resolution was electronically enhanced (voxel size 53.86 μ m) by Datos|x-reconstruction software (version 1.5.0., GE Sensing Inspection Technologies). Reconstructions (Supplementary Videos 1 and 2) were by Avizo (version 8.0, Visualization Science Group).

Two phylogenetic analyses were performed using PAUP (4.0b10) to estimate the placement of *Spinolestes* among eutriconodont mammals and their immediate outgroups, and among major cynodont-mammal clades. The analysis of

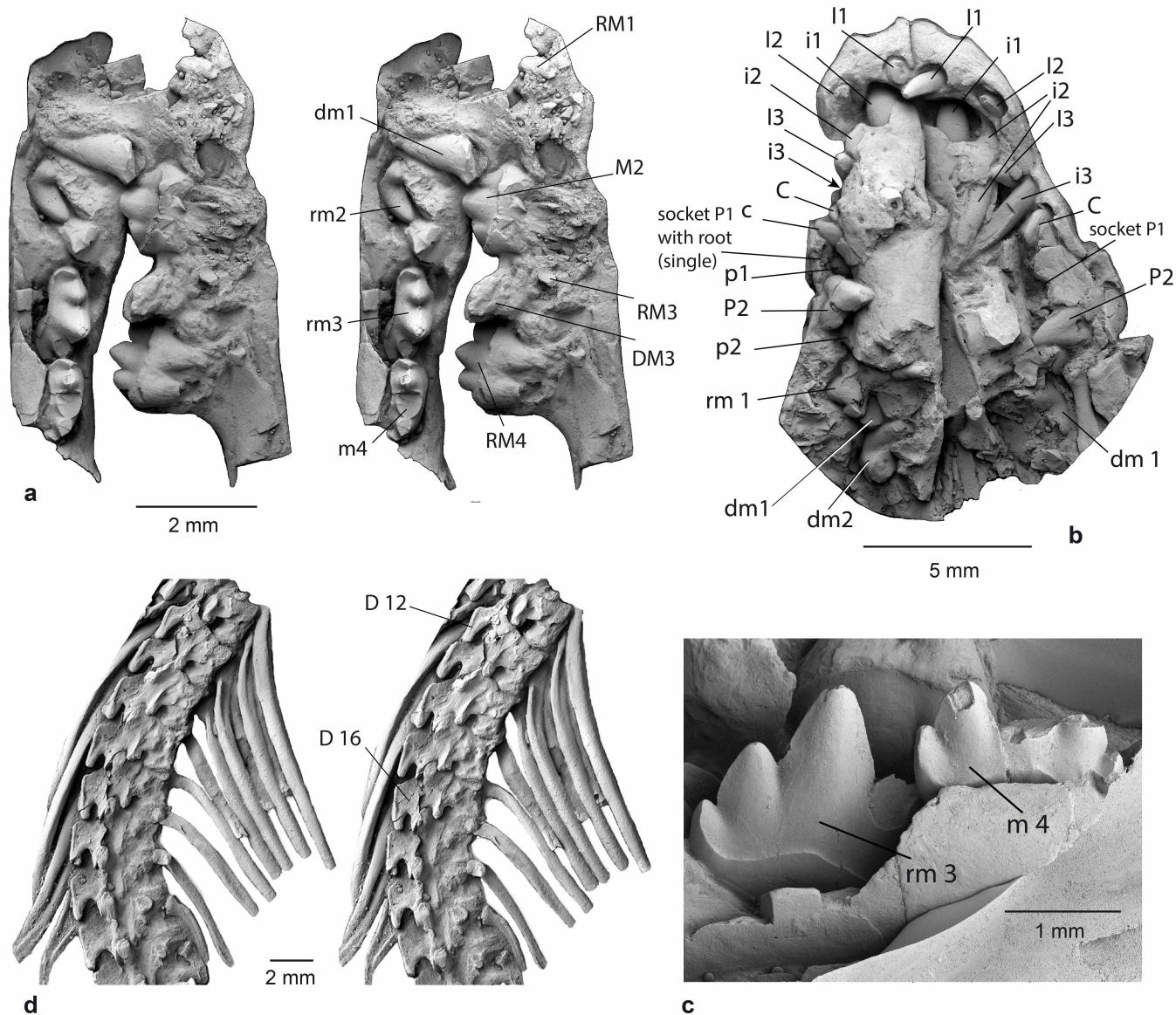
eutriconodonts was based on a data set augmented from refs 19, 20, 32, 33. The matrix has 29 taxa (20 eutriconodonts plus 9 outgroups) and 87 parsimony-informative characters. In this matrix, all characters have equal weight, six multi-state characters are ordered, and remaining characters are unordered. The branch-and-bound, which guarantees the shortest tree, yielded a single most parsimonious tree (tree length = 249; consistency index = 0.446; retention index = 0.755), as presented in Extended Data Fig. 3a. Nodal strengths of key clades are also assessed by bootstrap values (on the basis of a 50% majority rule consensus from 1,000 duplicates of bootstrap searches), and by Bremer index. Details in Supplementary Information.

Phylogenetic analysis of *Spinolestes* among major clades of cynodonts and mammals was based on a data set expanded and modified from that of ref. 7. The data matrix has 112 taxa (including 8 eutriconodonts), 490 characters (all parsimony-informative). The phylogenies were estimated by 1,000 replicates of heuristic search with a tree bisection and reconnection branch-swapping algorithm. All multi-state characters are unordered; all characters have equal weight. The search yielded 172 equally parsimonious trees (each tree has a length of 2,386 steps). A strict consensus tree of the 183 equally parsimonious trees is presented in Supplementary Information and a simplified tree is presented in Extended Data Fig. 3b.



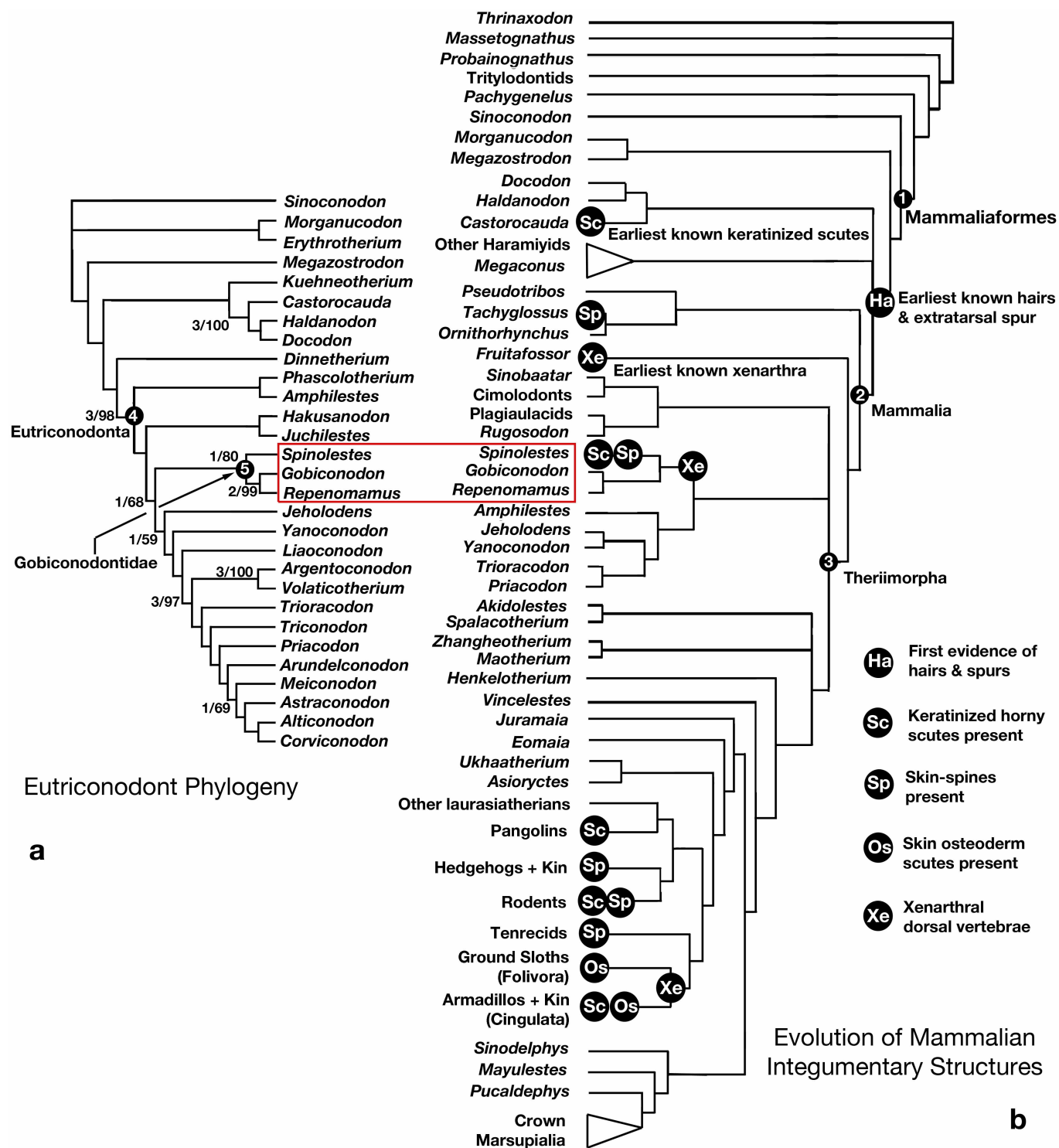
Extended Data Figure 1 | New gobiconodontid *S. xenarthrosus*, holotype transferred slab MCCMLH30000A. Skull exposed in ventral aspect. Inset: left calcaneus in dorsal view on the same slab somewhat apart from the skeleton.

Arrow points to skin patch preserving hair bundles between dorsal vertebrae 11 and 14.



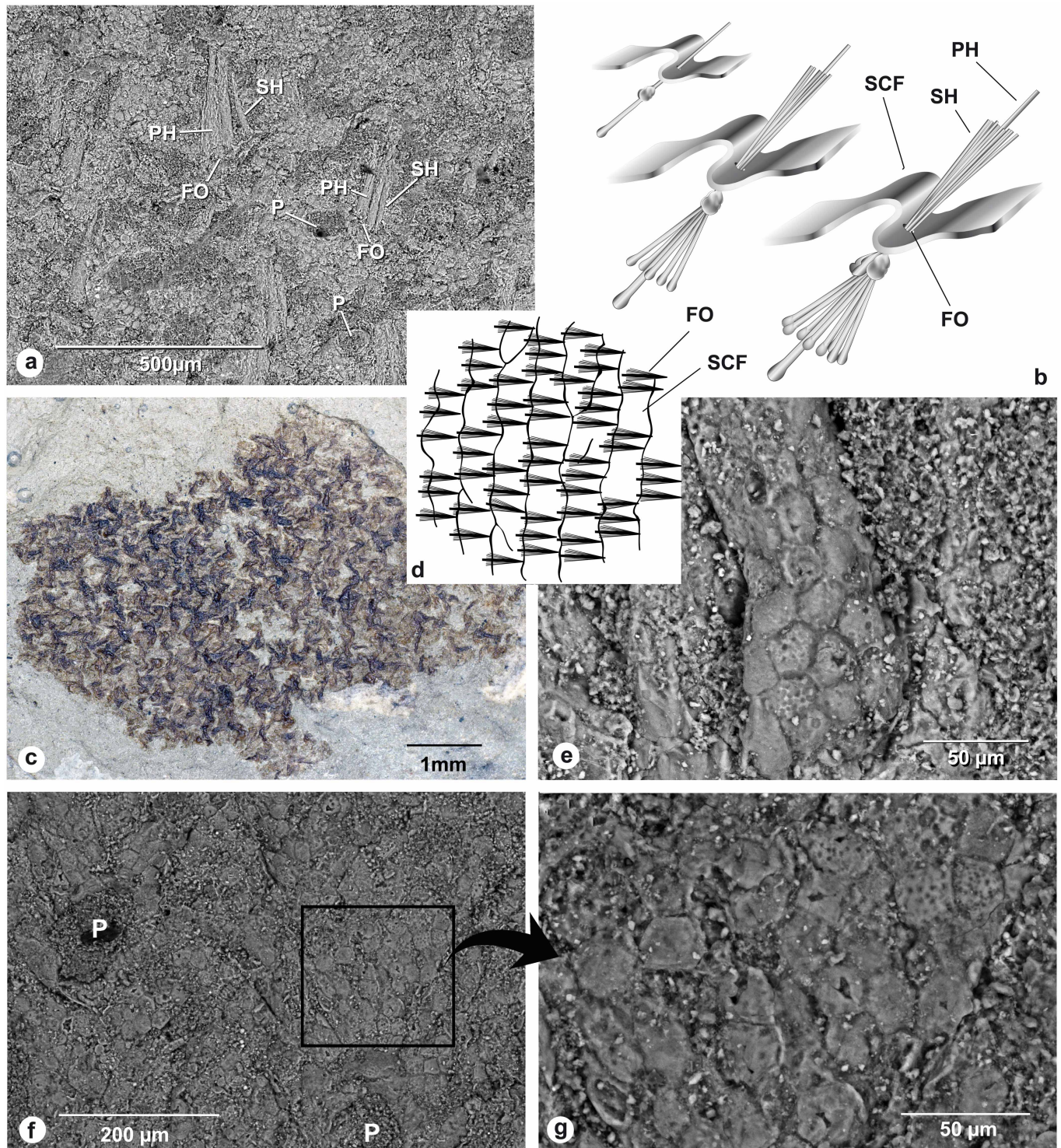
Extended Data Figure 2 | Dentition and xenarthrous vertebrae of *Spinolestes*. **a**, SEM images (stereo pair) of right lower mandible and right upper maxilla with molariforms. **b**, SEM image of anterior portion of skull and mandibles with dentition. **c**, SEM image of penultimate (m3) and ultimate (m4)

right lower molariforms in lingual view. Abbreviations: C/c, upper/lower canine; I/i, incisor; M/m, molariform; P/p, premolar; D/d, deciduous tooth; R/r, replacing tooth. **d**, SEM images (stereo pair) of dorsal vertebrae D 12 to D 19 in lateral aspect with ribs.



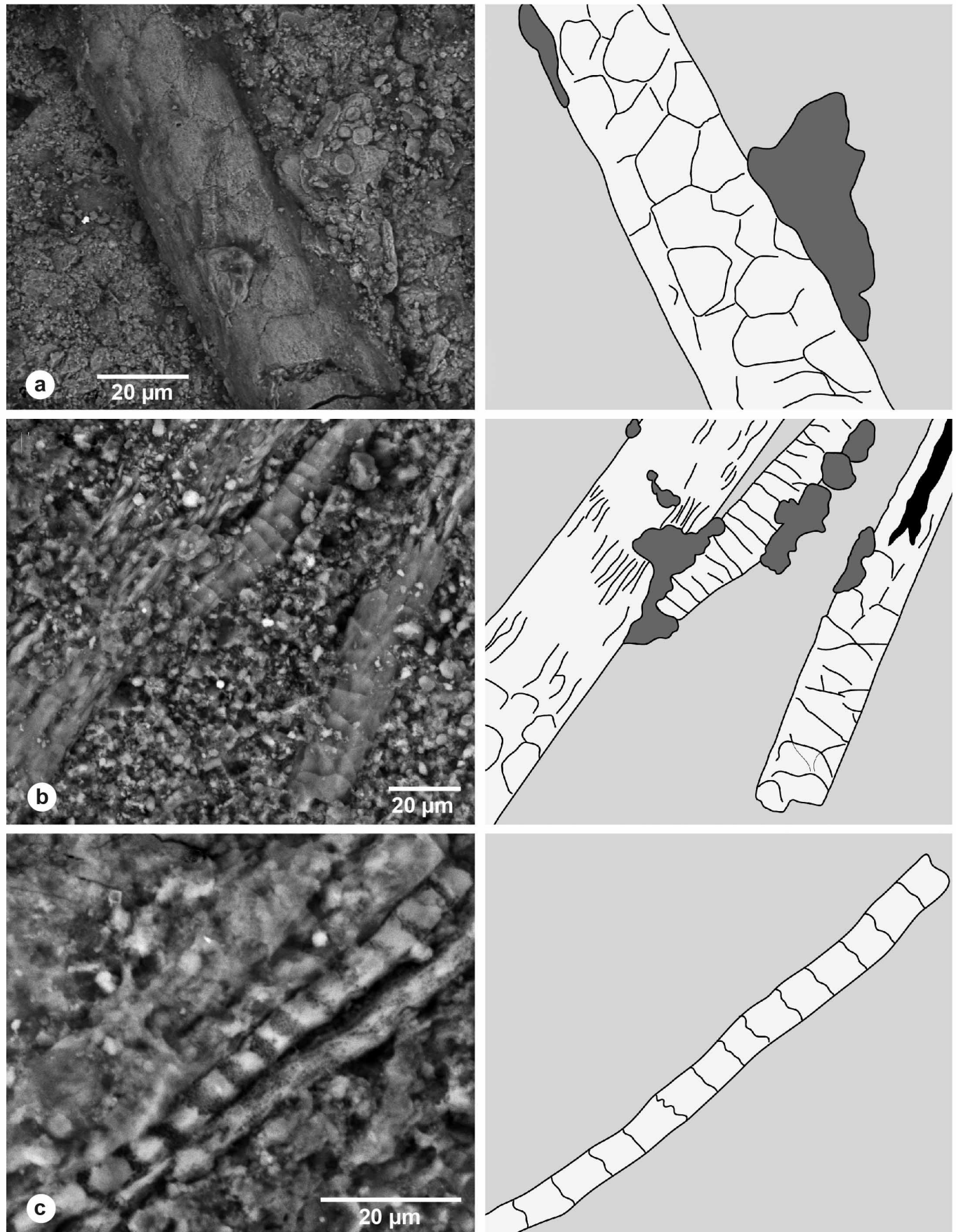
Extended Data Figure 3 | Phylogenetic relationships of *Spinolestes* and patterns of mammalian integumentary structure in early mammalian evolution. **a**, Position of *Spinolestes* within Eutriconodonta (Bremer index/

bootstrap values on the key clades of this study). **b**, Evolution of mammalian integumentary structures (simplified tree). Data sets and full analyses are presented in Supplementary Information.

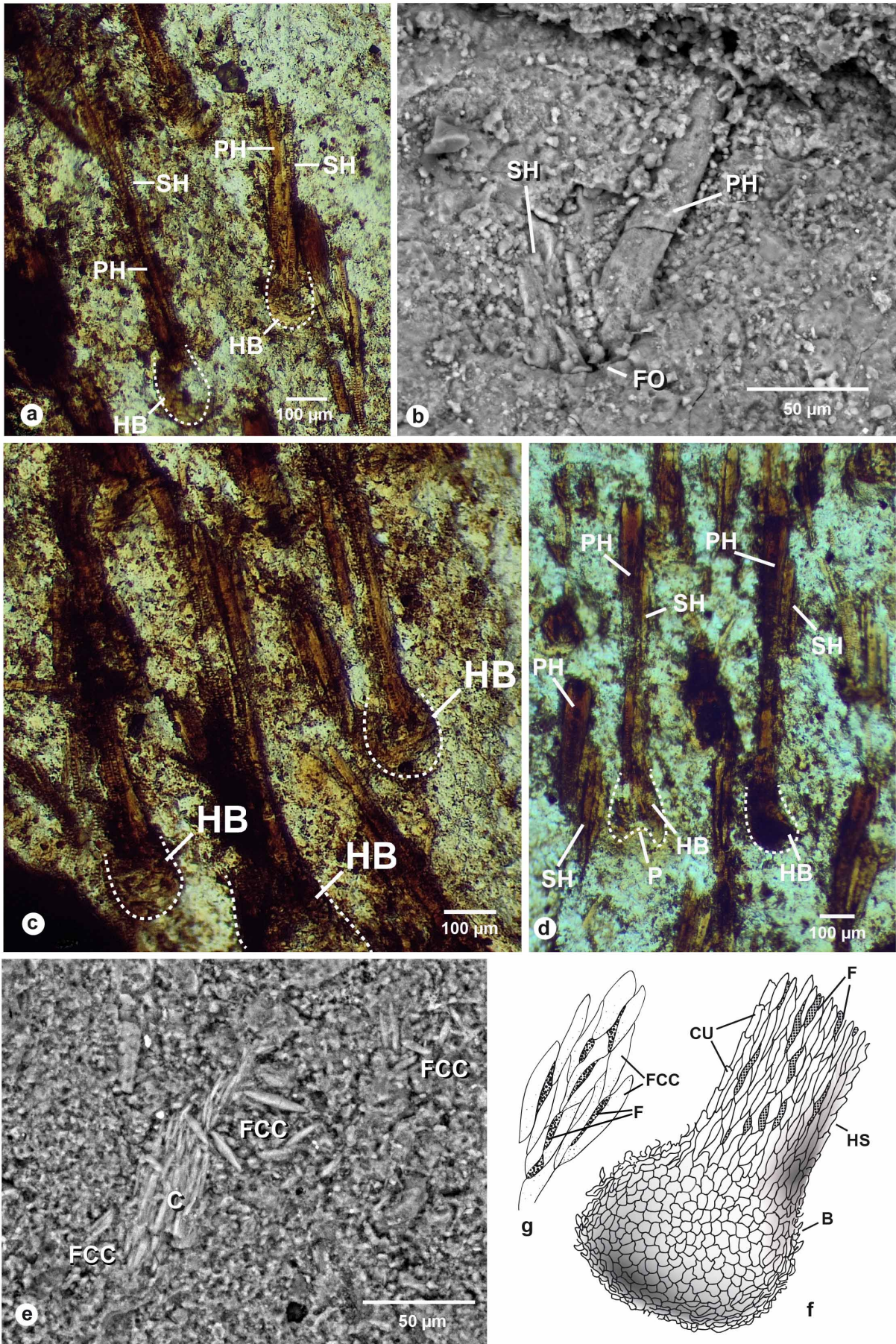


Extended Data Figure 4 | Skin structures of *S. xenarthrosus*. a, SEM image of skin surface showing compound hair follicles (FO), epidermal cells (keratinocytes), and pores (P). b, Schematic drawings of scale-like skin folds (SCF) with compound hair follicles (FO) of the dog with primary (PH) and secondary (SH) hairs of three ontogenetic stages (redrawn from ref. 38). c, Skin

surface of *Spinolestes* with scale-like wrinkles. d, Schematic diagram of scale-like wrinkled skin folds with hair follicles of a dog (redrawn from ref. 38). e, SEM image of epidermal cells (keratinocytes) of a hair follicle. f, SEM images of skin surface with polygonal epidermal cells (keratinocytes) and pores. g, Detail of the keratinocytes.

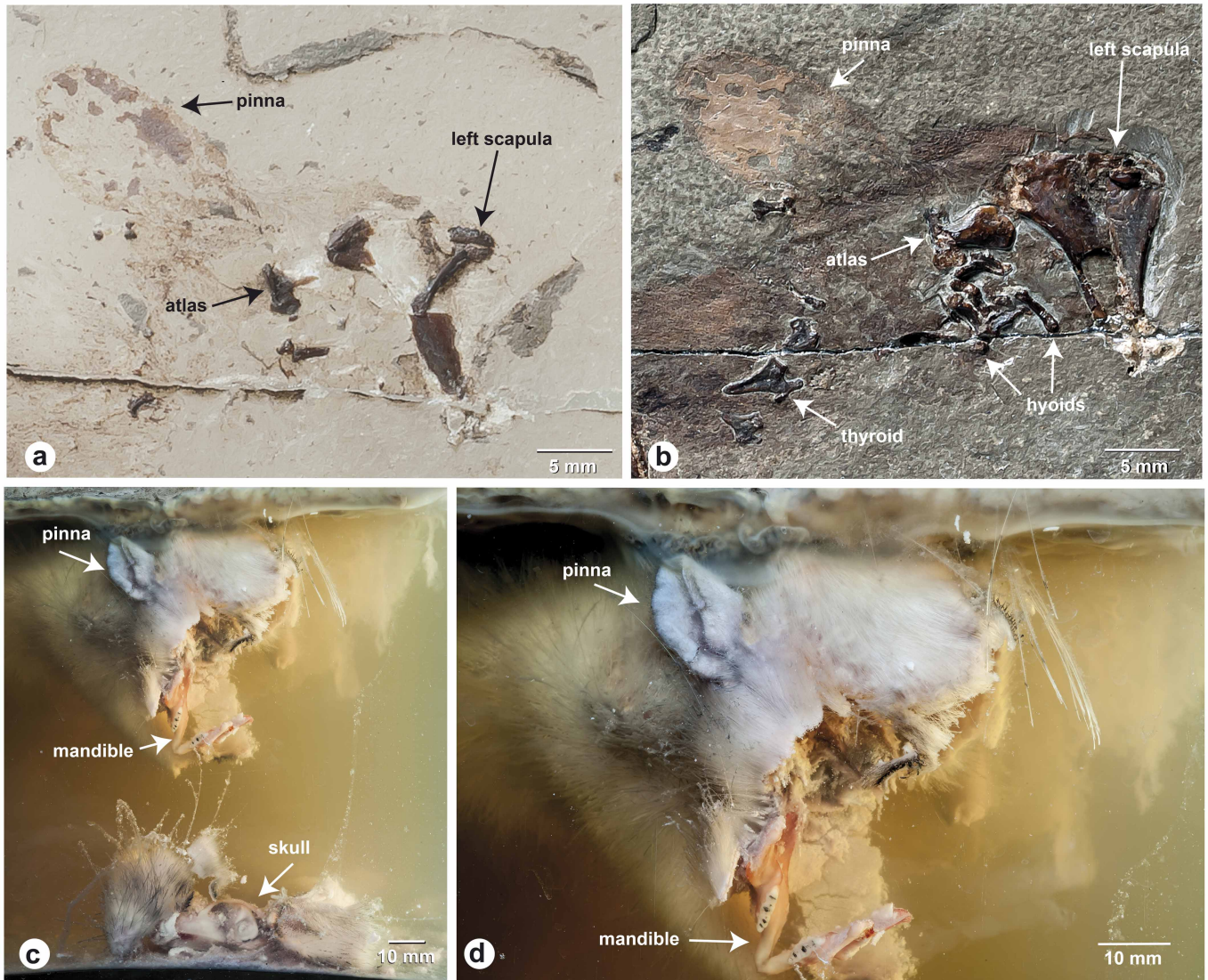


Extended Data Figure 5 | SEM images and interpretative line drawings of cuticular scale patterns. a–c, Primary (a, b) and secondary (b, c) hairs.



Extended Data Figure 6 | Hair structures of *S. xenarthrosus*. **a, c, d**, Patch of skin located between dorsal vertebrae 11 and 14 (arrow in Fig. 1 and Extended Data Fig. 1) on slab MCCMLH30000A under translucent light. **b**, SEM image of an orifice of a compound follicle (FO) with primary hair (PH) and three broken secondary hairs (SH). **e**, SEM image of a fraying hair of *S.*

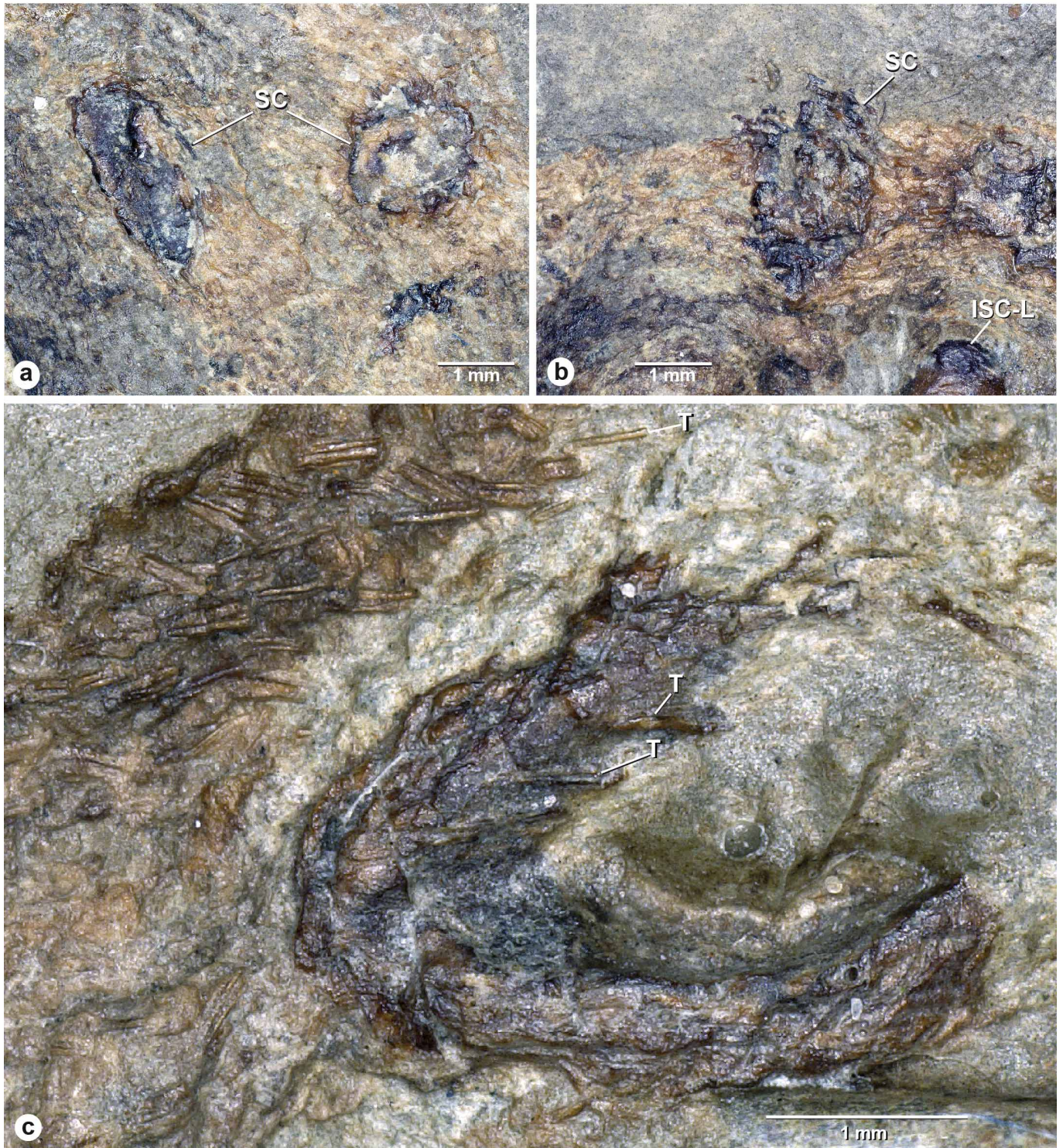
xenarthrosus showing fusiform cortical cells (FCC). **f**, Schematic diagram of human hair bulb (B) with cuticular scales (CU) and fusi (F) of the hair shaft (HS). **g**, Schematic diagram of a human head hair with fusiform cortical cells (FCC) and fusi. **f, g**, Redrawn from ref. 40. Abbreviations: C, cuticula; HB, hair bulb.



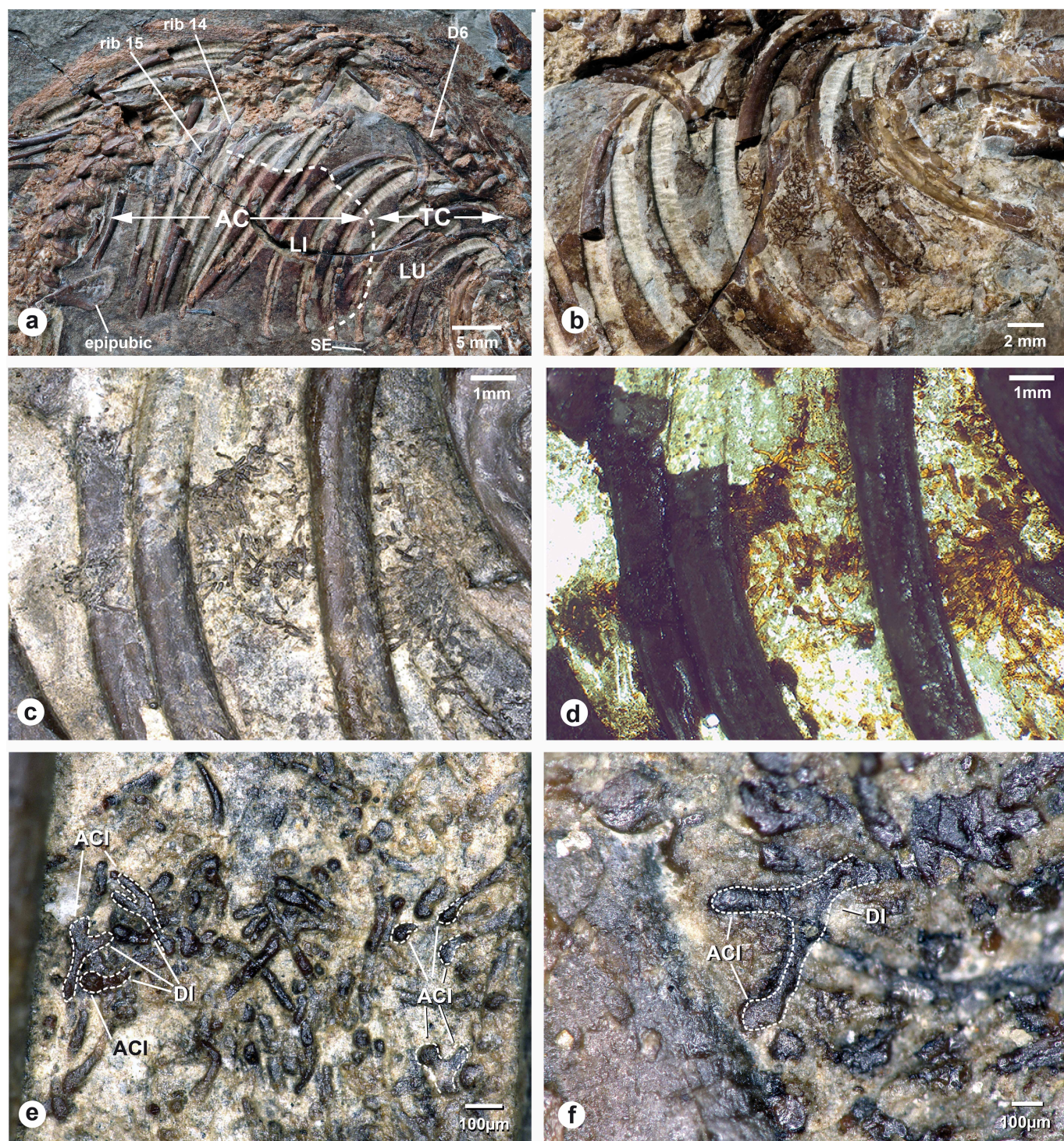
Extended Data Figure 7 | External pinna of *S. xenarthrosus*.

a, b, Comparison with pinna and scalp of decaying rat (*Rattus norvegicus*): **a**, on transferred slab MCCMLH30000A; **b**, on counter slab MCCMLH30000B on original limestone rock matrix. **c, d**, Pinna and scalp of decaying rat (*Rattus norvegicus*) for comparison. **c**, Head of decaying rat after 274 days in water at

room temperature. The skull is detached from the scalp and has fallen to the bottom. Scalp and right pinna are still intact and connected to the torso of the floating carcass. **d**, Detail of **c** with well-preserved right pinna in original position at the scalp. The mandible is displaced and loosely hanging down from the skin.



Extended Data Figure 8 | Keratinous dermal scutes (SC) of *S. xenarthrosus*. a, Oval dermal scutes dorsally of dorsal vertebra 20. b, Oval dermal scute dorsally of left ischium (ISC-L). c, Oval dermal scute from dorsal region with tubules (T) merging to homogeneous keratinous matrix.



Extended Data Figure 9 | Visceral cavities of *S. xenarthrosus* with internal organs. **a**, Detail of the visceral cavities of MCCMLH30000B with separation of the anterior thoracic cavity (TC) containing the lungs (LU), and the posterior abdominal cavity (AC) containing the liver (LI). Dashed line represents the diaphragm. SE, sternal elements. **b**, Lung tissue in the anterior part of the

thoracic cavity. **c**, Peripheral conducting and acinar airways of the lungs with typical dichotomous structure of the bronchial tree. **d**, Peripheral conducting and acinar airways of the lungs from transferred slab MCCMLH30000A under translucent light. **e**, **f**, Details of the peripheral conducting airway dichotomies (DI) and acinar structures (ACI) of the lungs.



Extended Data Figure 10 | Lifelike reconstruction of *S. xenarthrosus*.

Glia-derived neurons are required for sex-specific learning in *C. elegans*

Michele Sammut¹, Steven J. Cook², Ken C. Q. Nguyen², Terry Felton¹, David H. Hall², Scott W. Emmons^{2,3}, Richard J. Poole^{1*} & Arantza Barrios^{1*}

Sex differences in behaviour extend to cognitive-like processes such as learning, but the underlying dimorphisms in neural circuit development and organization that generate these behavioural differences are largely unknown. Here we define at the single-cell level—from development, through neural circuit connectivity, to function—the neural basis of a sex-specific learning in the nematode *Caenorhabditis elegans*. We show that sexual conditioning, a form of associative learning, requires a pair of male-specific interneurons whose progenitors are fully differentiated glia. These neurons are generated during sexual maturation and incorporated into pre-existing sex-shared circuits to couple chemotactic responses to reproductive priorities. Our findings reveal a general role for glia as neural progenitors across metazoan taxa and demonstrate that the addition of sex-specific neuron types to brain circuits during sexual maturation is an important mechanism for the generation of sexually dimorphic plasticity in learning.

During sexual maturation, the nervous system undergoes sexually dimorphic changes that couple behaviour to reproductive needs^{1–5}. Sex differences in behaviour extend beyond courtship and mating to cognitive-like processes such as learning that also enhance reproductive success^{6,7}. However, the precise dimorphisms in neural circuit development and organization that generate these differences in behavioural plasticity are mainly unknown.

C. elegans is a sexually dimorphic species with hermaphrodites that can self-fertilize and males that can only cross-fertilize hermaphrodites. As a consequence, males display a repertoire of reproductive behaviours not performed by the hermaphrodite and these behaviours require sex differences in the nervous system that develop during sexual maturation^{8,9}. The best described dimorphisms are the male-specific sensory and motor circuits required for mating. However, non-reproductive behaviours such as learning are also sexually dimorphic. Males, but not hermaphrodites, are capable of sexual conditioning, a process that involves associative learning and leads to a switch in chemotactic responses that facilitates effective mate finding¹⁰. Sexual conditioning requires sex-shared and male-specific sensory neurons¹⁰, but whether the circuit also contains sexually dimorphic interneurons for integration is not known.

The MCMs are a novel class of interneurons

We identified a previously unnoticed bilateral pair of neurons in the head of males that are required for sexual conditioning. We have termed them the MCMs (for mystery cells of the male) (Fig. 1). The identification of the MCMs is surprising given that the *C. elegans* anatomy, neuronal connectivity and developmental lineage have been characterized extensively^{8,9,11–14}. We noted the presence of the MCMs when we were analysing the expression pattern of a reporter transgene for the neuropeptide *pdf-1*. In males, but not in hermaphrodites, we observed *pdf-1* expression in a bilateral pair of cells located dorso-anterior to the pharyngeal metacarpus, a region devoid of neuronal cell bodies in hermaphrodites (Fig. 1a–c). This region is comprised of glial and epithelial cells, some of which project anteriorly to the

nose^{13,15}. Unlike any other cells in this region, we observed that the male's *pdf-1*-positive cells send projections posteriorly into the nerve ring (Fig. 1b and Extended Data Fig. 1a, b). The nerve ring is the main concentrated neuropil of the nematode brain, where most sensory neuron–interneuron synapses and neural processing occur¹³. We found that the *pdf-1*-positive cells also express the pan-neuronal reporter *rab-3*, suggesting they may be neurons. *rab-3* expression was first observed at the L4 stage, when sexual maturation begins (Fig. 1c). Through reporter gene analysis, and ultrastructural reconstruction, we found that the MCMs are indeed fully differentiated neurons. In addition to *rab-3*, these cells express a battery of neuronal genes required for both electrical and chemical communication (Fig. 1d, Extended Data Fig. 1c and Extended Data Table 1). These include components of the SNARE complex; Na⁺ and voltage-gated Ca²⁺ channel subunits; innexins; components of the machinery for neuropeptide secretion *ric-19* and *ida-1*; and *pdf-1*. At the ultrastructural level, we observed dense-core vesicles in the MCM somata and projections, and clear-core vesicles at presynaptic densities (Fig. 1e). Together, these observations indicate that the MCMs are a novel class of male-specific interneurons.

We identified the MCMs in electron micrographs of serial sections through the male head and established their entire neural connectivity (Fig. 2 and Extended Data Table 2). The cell bodies send processes posteriorly in the amphid bundles, from which they diverge to enter the nerve ring, and from which they exit and extend into the ventral nerve cord (Extended Data Fig. 1b). In the nerve ring and ventral nerve cord we identified synaptic interactions with 24 neuron classes (Extended Data Fig. 1d and Extended Data Table 2). The bulk of the synaptic input (65%) is from three interneuron classes, two of which are sex-shared (AVF and PVQ) and one of which is male-specific (EF). These interneurons receive extensive sensory inputs from the male copulatory circuits in the tail and extend processes through the ventral nerve cord and into the nerve ring, where they connect to the MCMs both directly and through RIF (a sex-shared class of second-order interneurons that receive input from head

¹Department of Cell and Developmental Biology, University College London, London WC1E 6BT, UK. ²Dominick P. Purpura Department of Neuroscience, Albert Einstein College of Medicine, Bronx, New York 10461, USA. ³Department of Genetics, Albert Einstein College of Medicine, Bronx, New York 10461, USA.

*These authors contributed equally to this work.

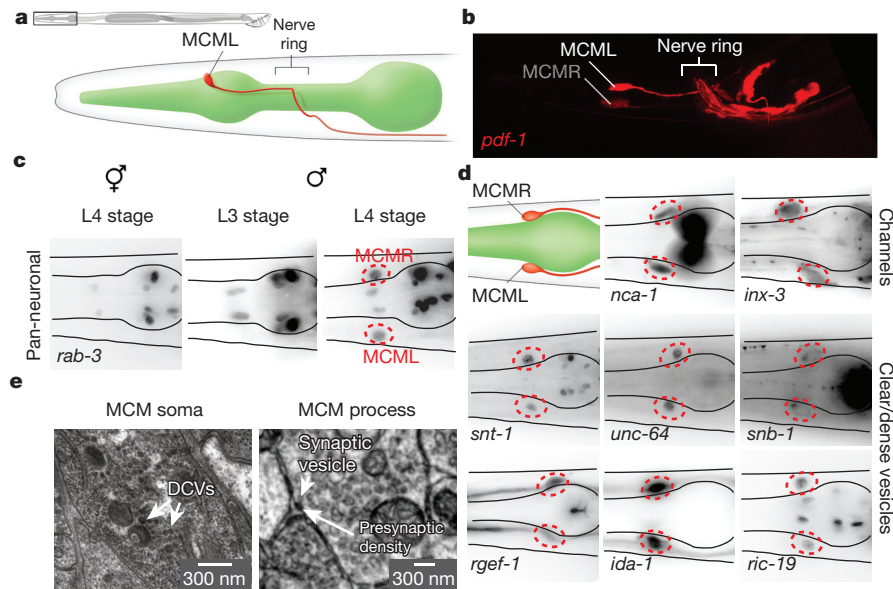


Figure 1 | The MCMs are newly identified male-specific neurons. **a–d**, Lateral views (**a**, **b**) and dorsal views (**c**, **d**) of animals oriented anterior to the left. **a**, WormAtlas-style⁴⁷ diagram depicting the morphology and position, adjacent to the pharynx, of one of the bilateral pair of MCM neurons in the head of a male. **b**, Confocal projection of *pdf-1::rfp* expression in the head of an adult male (same region as in **a**). **c**, Expression of the pan-neuronal reporter transgene *rab-3::rfp* (Ras GTPase) in the head of a hermaphrodite and males at the third (L3) and fourth (L4) larval stages. The position of the

MCMs is indicated with dashed red circles. L, left; R, right. **d**, Expression in the MCMs in adult males of reporter transgenes for neuronal markers. *nca-1* (NALCN Na⁺ channel subunit); *inx-3* (gap junction innexin); *snt-1* (synaptotagmin); *unc-64* (syntaxin); *snb-1* (synaptobrevin); *rgef-1* (Ras exchange factor); *ida-1* (tyrosine phosphatase-like receptor, phogrin); *ric-19* (cytosolic, vesicle secretion). **e**, Electron micrographs of MCM showing dense-core vesicles (DCVs) (left) and a synapse (right).

chemosensory circuits) (Fig. 2). The largest synaptic output is onto the AVB pre-motor interneurons that drive forward locomotion. These outputs are both direct (41%) and indirect via RIF (14%). Thus, the presence of the MCMs creates a series of male-specific disynaptic feed-forward triplet motifs connecting male-specific EF and sex-shared AVF, PVQ and RIF to AVB (Fig. 2). Such triplet

motifs recur frequently in the *C. elegans* nervous system^{9,16}. There is also a small amount of reciprocal interaction between the MCMs and the male-specific, pheromone-sensing CEM head sensory neurons. The major input from mate-sensing circuits and their output onto second- and third-order interneurons make the MCMs ideally placed for the integration of mate-experience into circuits regulating behavioural plasticity to sensory stimuli.

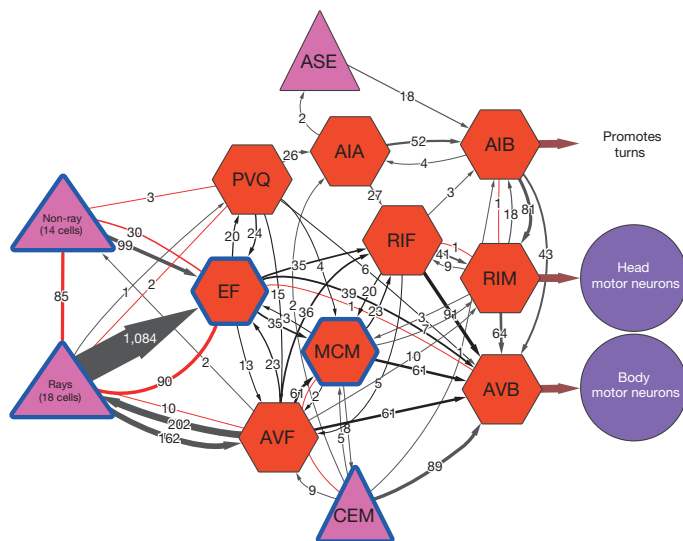


Figure 2 | MCM connectivity. **a**, Connectivity diagram of the MCMs showing the main inputs and outputs. The connections to the neurons known to regulate sexual conditioning (rays and CEMs), chemosensory plasticity (AIA, RIF), and salt sensation (ASE) are included. Grey and red connections indicate chemical and electrical synapses, respectively. The triplet motifs created by the MCMs are highlighted in black. The thickness of the lines is proportional to the anatomical strength of the connections, indicated in numbers (the number of EM serial sections of scored synaptic connectivity). Pink triangles, sensory neurons; red octagons, interneurons; purple circles, motor neurons; blue outline, male-specific.

The MCMs regulate sex-specific learning

We found that the MCMs are required for a male-specific switch in chemosensory behaviour induced by sexual conditioning. Sexual conditioning is a mate-experience-dependent process of behavioural modification that overrides the effects of starvation during chemosensory learning¹⁰. This plasticity confers on males the ability to use recent environmental chemosensory cues as signals for effective mate finding, and reflects the adult male prioritization of sex over food^{17,18}. One chemosensory behaviour that is subject to sexual conditioning is salt-avoidance learning¹⁰. *C. elegans* is normally attracted to salt but both males and hermaphrodites can learn to avoid salt when it is previously associated with an aversive stimulus such as starvation^{19,20} (Fig. 3a). However, unlike hermaphrodites, males switch their behaviour and become attracted to salt if mates are present during previous conditioning¹⁰ (Fig. 3a). Ablation of the MCMs at the late fourth larval (L4) stage with a laser microbeam resulted in males that failed to undergo a sexually conditioned switch, thus avoiding salt rather than being attracted to it, after conditioning with starvation, salt and mates (Fig. 3b). Importantly, MCM ablation did not disrupt salt attraction in non-conditioned males or salt avoidance learning after conditioning with salt and starvation in the absence of mates (Fig. 3b), indicating that the defects are specific to sexual conditioning and not due to inability to sense salt or starvation.

Sexual conditioning requires two types of mate inputs¹⁰: a secreted pheromone synthesized by the enzyme DAF-22 and sensed through chemosensory amphid neurons and the male-specific CEM sensory neurons in the head^{10,21}, and a contact-dependent cue that is sensed through the male-specific ray neurons in the tail^{10,22}. MCM-ablated

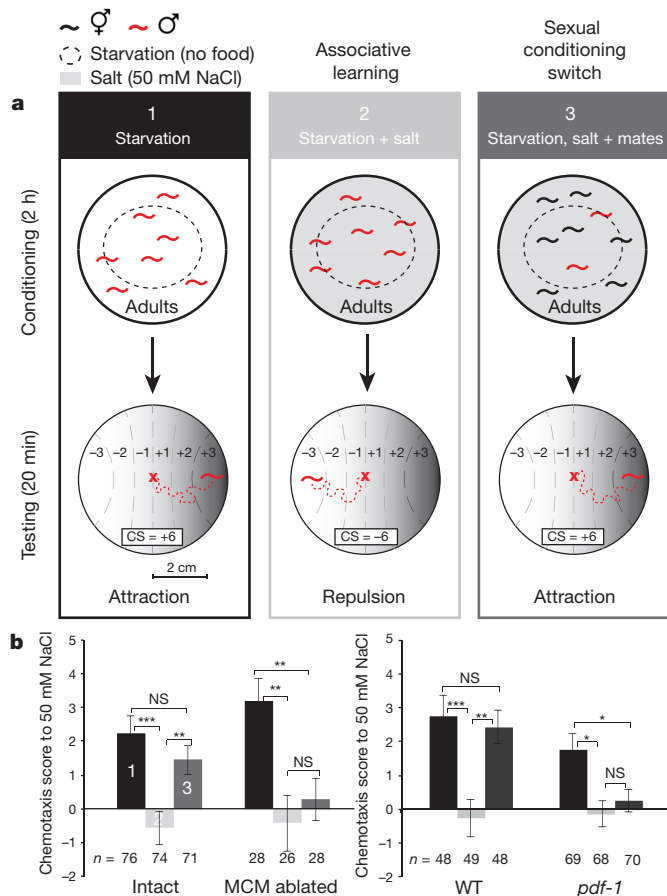


Figure 3 | The MCMs are required for male-specific associative learning. **a**, Diagram depicts the salt-avoidance learning and sexual conditioning assay. CS, chemotaxis score. **b**, Salt chemotaxis scores of previously conditioned intact and MCM-ablated males [*inIs179(ida-1::gfp);him-8(e1489)*], and wild-type [*him-5(e1490)*] and *pdf-1(tm1996)* mutant males. *n*, number of individual animals tested. Error bars indicate s.e.m. Mann-Whitney *U*-test was used for statistical analysis. ****P* < 0.001; ***P* < 0.01; **P* < 0.05; NS = no statistically significant difference (*P* ≥ 0.05). WT, wild type.

males responded efficiently to both *daf-22*-dependent secreted pheromones (ascarosides) (Fig. 4a, b and Extended Data Fig. 2a) and contact-dependent hermaphrodite cues that initiate the mating sequence and inhibit male food-leaving behaviour^{22,23} (Fig. 4c, d and Extended Data Fig. 2b, f). These results indicate that the MCMs are not required for the sensory detection of mates, salt or starvation, but that they specifically regulate the integration of these cues to produce behavioural plasticity in a context- and sex-specific manner.

Having found a role for the MCMs in sexual conditioning, we next asked whether the MCMs are important regulators of other male-specific behaviours from the coordinated sensory-motor programmes of mating²⁴ to other integrative behaviours such as the decision to leave food in search of mates^{22,25}. MCM ablation did not cause defects in any other male-specific behaviour tested or overall fertility (Fig. 4c, d and Extended Data Fig. 2b–f), demonstrating a surprisingly specific role for the MCMs in sexual conditioning.

To establish whether sexual conditioning is regulated through neuromodulation, we tested *pdf-1(tm1996)* mutants in the salt chemotaxis learning assay. The *tm1996* deletion allele removes the minimal promoter and first exon of the *pdf-1* gene and is likely to be a null^{26,27}. We found that *pdf-1* mutant males, like MCM-ablated males, failed to undergo sexual conditioning and avoided salt after conditioning with starvation, salt and mates (Fig. 3b). Loss of *pdf-1*, like loss of MCMs, did not disrupt salt sensation in naive animals or the ability to associate salt and starvation (Fig. 3b). *pdf-1* mutants also responded to secreted and contact-dependent mate cues, although they were slightly, but significantly, less efficient than wild-type males at contact response (Fig. 4a, b and Extended Data Fig. 2b)²⁷. Together, these results tentatively suggest that the MCMs may modulate the circuits that regulate chemosensory plasticity through neuropeptide secretion.

The MCMs arise from glia

Next, we sought to establish the developmental mechanisms by which these male-specific neurons arise. Unlike any other neurons in *C. elegans* and other invertebrates, which arise from epithelial or undifferentiated blast cells^{8,11,12,28}, we found that the MCMs arise from glial cells. Glia are specialized cells with projections that ensheath the cilia of the sensory neurons with which they are associated, and provide structural and functional support to these neurons²⁹. As neuronal gene expression in the MCMs begins during sexual maturation, we

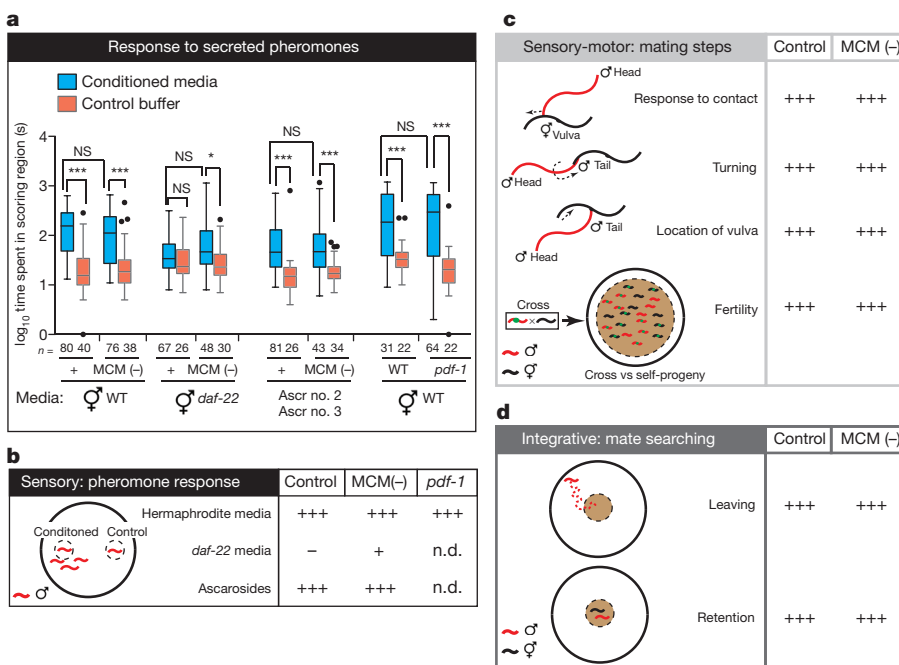


Figure 4 | MCM ablation does not affect other male-specific behaviours. Efficiency of intact and MCM-ablated males [*inIs179(ida-1::gfp);him-8(e1489)* or *otIs356(rab-3::rfp);him-5(e1490)*], and wild-type [*him-5(e1490)*] and *pdf-1(tm1996)* mutant males in male-specific behaviours.

a, Response to wild-type or *daf-22* (*m130*) hermaphrodite-conditioned media and purified ascaroside pheromones (80 nM Ascr no. 3, 800 nM Ascr no. 2). Graphs represent Tukey box plots of logarithmic transformations of the data; *n* = number of independent events (that is, entry in scoring region); *t*-test with Bonferroni correction was used for statistical analysis. ****P* < 0.001; **P* < 0.05; NS, no statistically significant difference (*P* ≥ 0.05). **b**, Summary diagram of the data plotted in **a**. Statistical significance compared to control buffer is indicated by + + +, *P* < 0.001; +, *P* < 0.05; -, no statistically significant difference. n.d., not determined. **c**, **d**, Execution of mating sub-steps (**c**) and exploration in search of mates (**d**). + + + indicates performance level of intact control males.

examined the expression of an *rnr-1::gfp* reporter, which labels cells in S-phase³⁰, to determine whether the MCMs are born at this stage via an undescribed cell division in the male head. In the head region, we observed strong expression only in males (Fig. 5a). Expression was observed in a bilateral pair of cells at the late L3 stage, corresponding to the amphid socket (AMso) glial cells, and perduring in two bilateral pairs of cells at the early L4 stage, corresponding to the two AMso and the MCMs (Fig. 5b and Extended Data Fig. 3a). The lack of expression in the hermaphrodite head is consistent with the lack of AMso glial cell division in this sex¹².

Ablation of one of the bilateral pair of AMso glial cells at the L3 stage in males resulted in a consistent loss of the MCM on the operated side (Extended Data Table 3). Importantly, the MCM was not lost when the AMso was ablated at the mid-L4 stage, after the MCM was born (Extended Data Table 3). Together this demonstrates that the AMso glial cells are the MCM progenitors.

The AMso glial cells display a polarized morphology with a single projection running anteriorly along the anterior to posterior (A-P)

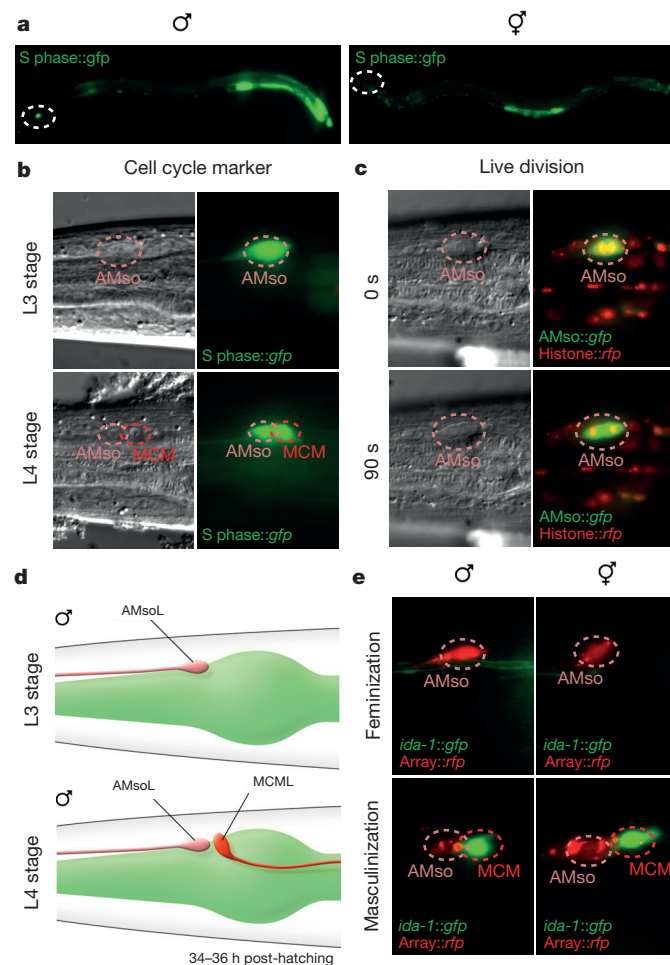


Figure 5 | The MCMs originate from a male-specific cell division of the AMso glial cells. Lateral views of animals oriented anterior to the left and dorsal to the top. **a**, **b**, Expression of the S phase reporter transgene *rnr-1::gfp* in whole animals at the L4 stage (**a**) and in the head of L3 and L4 males (**b**). The male-specific division in the head is indicated by dashed circles in **a**. The rest of expression corresponds to developing reproductive structures. **c**, Differential interference contrast (DIC) and fluorescent images of the AMso cell body with RFP-labelled histones during cell division. **d**, WormAtlas diagram depicting the morphology and position of the AMso and the MCM in the male head at the L3 and L4 stages. **e**, Head of animals with sex-reversal genetic manipulations of AMso. Feminization by expression of the *grl-2::tra-2(IC)::SL2::mCherry* transgene *oleEx23*. Masculinization by expression of *grl-2::fem-3::SL2::mCherry* transgenes *oleEx18* and *oleEx24*.

axis to the nose¹⁵ (Figs 5d and 6a, b). Live imaging of the AMso cell division in males revealed that the cleavage plane was perpendicular to the axis of AMso polarity and that the AMso projection was maintained throughout the division (Fig. 5c and Extended Data Fig. 3b). The projection was inherited by the self-renewing AMso glial daughter.

AMso plasticity is intrinsically regulated

The AMso glial cells are born during embryogenesis from the same fully described lineages in both sexes (Extended Data Fig. 4a)¹¹, but retain their plasticity only in males, re-entering the cell cycle during sexual maturation. We found that the AMso developmental program occurred according to the genetic sex of the AMso cells and not according to that of the rest of their lineage or the animal. Sexual dimorphism in *C. elegans* is regulated by a genetic pathway that converges on the Ci/GLI-like zinc-finger transcription factor *tra-1*, which inhibits male development and activates hermaphrodite development^{31,32}. Consistent with a role for the sex-determination pathway in MCM specification, 22/22 *tra-1(e1488)* sex-transformed XX pseudo-males produced MCMs (data not shown). The sister cells of the AMso are the male-specific CEM sensory neurons, which undergo cell death during embryogenesis in hermaphrodites¹¹ (Extended Data Fig. 4a). In order to establish whether sex-specific neurogenic competence of the AMso glial cells is specified extrinsically

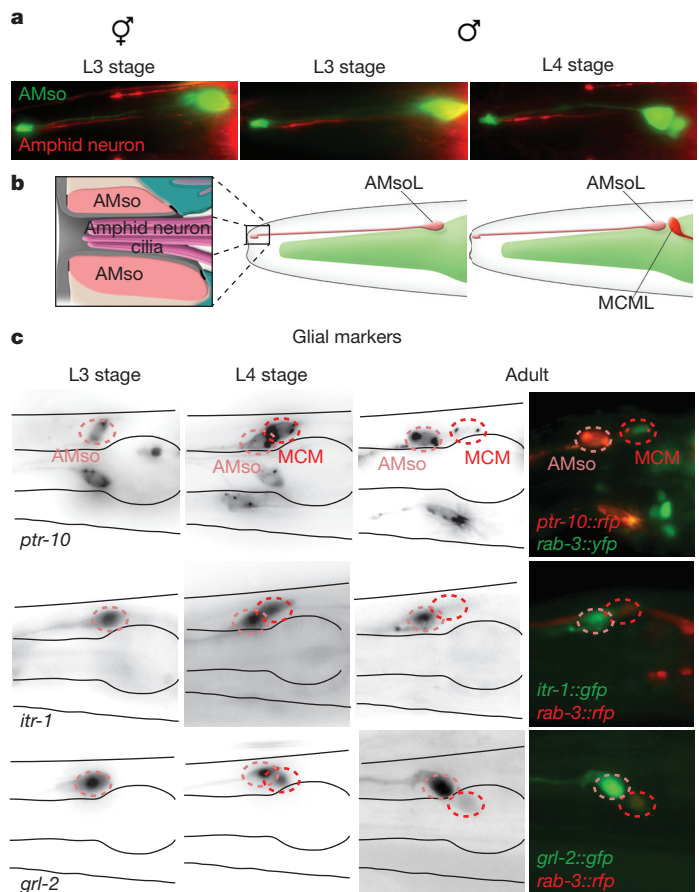


Figure 6 | The male AMso cells are fully differentiated glia before and after the division that generates the MCM neuron. Lateral views of animals oriented anterior to the left and dorsal to the top. **a**, AMso glial cell (green, *grl-2::gfp*) and amphid sensory neuron dendrites (red, *pdf-1::rfp*) in an L3 hermaphrodite and L3 and L4 males. **b**, Diagram of the AMso glial cell and the distal end of its projection to the nose, where it ensheaths the cilia of sensory dendrites. **c**, Expression in AMso and perdurance in MCM of reporter transgenes for glial markers in L3, L4 and adult males. *ptr-10* (Patch-related receptor); *itr-1* (IP3 receptor); *grl-2* (Hedgehog-like/Ground-related).

or intrinsically, we expressed sex-reversing transgenes under cell-restricted promoters to uncouple the genetic sex of the AMso glial cells from that of the rest of the animal. We used the *grl-2* and *ztf-16* promoters, which drive expression in the AMso cells after they are born and one or two other non-overlapping cell types in the head in larval animals^{33,34}. Masculinization of the AMso cells in hermaphrodites, by expression of a *fem-3* transgene^{35,36}, resulted in AMso cell division and MCM neuronal differentiation (Fig. 5e, Extended Data Fig. 4b and Extended Data Table 4). Conversely, feminization of AMso cells in males, by overexpression of the *tra-2* intracellular domain (*tra-2IC*)^{18,37}, resulted in lack of AMso cell division and lack of MCMs (Fig. 5e, Extended Data Fig. 4c and Extended Data Table 4). Masculinizing transgenes in males or feminizing transgenes in hermaphrodites had no effect on AMso development (Fig. 5e). These results indicate that the competence of the AMso glial cell to become a neural progenitor is specified intrinsically, and extend the battery of sexual dimorphisms regulated by the sex-determination pathway to the intrinsic properties of glia.

The male AMso cells are fully differentiated glia

We found that the male AMso glial cells are fully differentiated specialized cells before they divide, indicating that the production of the MCMs from AMso involves a complete switch in terminal cell fate (Fig. 6). In hermaphrodites and adult males, the AMso cells are specialized glial cells that form a hole in the cuticle through which most amphid neurons contact the outside world¹⁵, (Fig. 6a, b). Hermaphrodite AMso cells also express a battery of glial-cell markers³⁸. We observed that at the L3 stage, before their division, the male AMso cells display the same characteristics as the AMso of hermaphrodites and adult males: they project to the nose where they form a socket that ensheaths amphid neuron cilia (Fig. 6a, b), and they express AMso markers (Fig. 6c and Extended Data Table 5). Following AMso division, one daughter retains the AMso identity and the other becomes the MCM, which loses molecular and structural characteristics of glia (Extended Data Fig. 5) and, as described above, acquires a neuronal identity. Similar complete switches in terminal cell fates have been described as transdifferentiation events in *C. elegans* and other systems^{11,39–41}.

Discussion

The production of the MCMs from AMso glial cells through a division that results in self-renewal and neuronal differentiation is equivalent to a snapshot of the reiterative neuronal differentiation events arising from radial glia cell divisions during vertebrate neurogenesis⁴². However, the extent to which vertebrate neural progenitors are fully differentiated glia remains unclear and, in particular, it is not known whether the support and progenitor functions of glia can be integrated within one cell^{43,44}. Our results provide the first example of neurons arising from glia in a non-vertebrate organism and demonstrate that fully differentiated, functional glia can retain neural progenitor properties during normal development.

In *C. elegans* the best described transdifferentiation event is the direct conversion of the rectal epithelial cell Y into a neuron^{39,45}. Other cell fate switches in *C. elegans* and other systems, which are not direct but require cell division, have also been suggested to be transdifferentiation events^{11,40,41}. The complete cell fate switch we observe suggests that the production of neurons from glia may be a process of natural transdifferentiation. It will be interesting to determine whether direct and indirect transdifferentiations of this kind occur through similar molecular mechanisms.

Our findings provide a direct link between developmental and anatomical sexual dimorphism in higher-order processing areas of the brain and sexually dimorphic behavioural plasticity during learning. We have shown that sex differences in learning require the production, late in development, of a male-specific class of cephalic

interneurons that were not previously known to exist. These interneurons are incorporated into pre-existing, sex-shared circuits allowing the male to change behavioural priorities according to its new reproductive needs. The MCMs function to confer salience to previous mate-experience during male navigation. Based on their connectivity, we suggest that these interneurons may function in other chemosensory-plasticity behaviours that may be subject to sexual conditioning. The production of more neurons in males is also a mechanism to generate cognitive sexual dimorphism in the song-learning system of songbirds, where neural progenitors have a glial identity⁴⁶. However, whether these neurons are a male-specific class and whether their progenitors are intrinsically sexually dimorphic is not known. Our findings indicate that the addition of sex-specific neurons is an effective way of remodelling brain circuits during sexual maturation and reveal a general, possibly ancient role for glia as neural progenitors to assemble circuits for higher-order processing.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 17 May; accepted 10 September 2015.

- Kimura, K.-I., Ote, M., Tazawa, T. & Yamamoto, D. Fruitless specifies sexually dimorphic neural circuitry in the *Drosophila* brain. *Nature* **438**, 229–233 (2005).
- Ruta, V. et al. A dimorphic pheromone circuit in *Drosophila* from sensory input to descending output. *Nature* **468**, 686–690 (2010).
- Yang, C. F. et al. Sexually dimorphic neurons in the ventromedial hypothalamus govern mating in both sexes and aggression in males. *Cell* **153**, 896–909 (2013).
- Rideout, E. J., Dorman, A. J., Neville, M. C., Eadie, S. & Goodwin, S. F. Control of sexual differentiation and behavior by the doublesex gene in *Drosophila melanogaster*. *Nature Neurosci.* **13**, 458–466 (2010).
- Stowers, L. & Logan, D. W. Sexual dimorphism in olfactory signaling. *Curr. Opin. Neurobiol.* **20**, 770–775 (2010).
- Nottebohm, F. & Arnold, A. P. Sexual dimorphism in vocal control areas of the songbird brain. *Science* **194**, 211–213 (1976).
- Keleman, K. et al. Dopamine neurons modulate pheromone responses in *Drosophila* courtship learning. *Nature* **489**, 145–149 (2012).
- Sulston, J. E., Albertson, D. G. & Thomson, J. N. The *Caenorhabditis elegans* male: postembryonic development of nongonadal structures. *Dev. Biol.* **78**, 542–576 (1980).
- Jarrell, T. A. et al. The connectome of a decision-making neural network. *Science* **337**, 437–444 (2012).
- Sakai, N. et al. A sexually conditioned switch of chemosensory behavior in *C. elegans*. *PLoS ONE* **8**, e68676 (2013).
- Sulston, J. E., Schierenberg, E., White, J. G. & Thomson, J. N. The embryonic cell lineage of the nematode *Caenorhabditis elegans*. *Dev. Biol.* **100**, 64–119 (1983).
- Sulston, J. E. & Horvitz, H. R. Post-embryonic cell lineages of the nematode, *Caenorhabditis elegans*. *Dev. Biol.* **56**, 110–156 (1977).
- White, J. G., Southgate, E., Thomson, J. N. & Brenner, S. The structure of the nervous system of the nematode *Caenorhabditis elegans*. *Phil. Trans. R. Soc. Lond. B* **314**, 1–340 (1986).
- Hall, D. H. & Russell, R. L. The posterior nervous system of the nematode *Caenorhabditis elegans*: serial reconstruction of identified neurons and complete pattern of synaptic interactions. *J. Neurosci.* **11**, 1–22 (1991).
- Ward, S., Thomson, N., White, J. G. & Brenner, S. Electron microscopical reconstruction of the anterior sensory anatomy of the nematode *Caenorhabditis elegans*. *J. Comp. Neurol.* **160**, 313–337 (1975).
- Varshney, L. R., Chen, B. L., Paniagua, E., Hall, D. H. & Chklovskii, D. B. Structural properties of the *Caenorhabditis elegans* neuronal network. *PLoS Comput. Biol.* **7**, e1001066 (2011).
- Barrios, A. Exploratory decisions of the *Caenorhabditis elegans* male: a conflict of two drives. *Semin. Cell Dev. Biol.* **33**, 10–17 (2014).
- Ryan, D. A. et al. Sex, age, and hunger regulate behavioral prioritization through dynamic modulation of chemoreceptor expression. *Curr. Biol.* **24**, 2509–2517 (2014).
- Saeki, S., Yamamoto, M. & Iino, Y. Plasticity of chemotaxis revealed by paired presentation of a chemoattractant and starvation in the nematode *Caenorhabditis elegans*. *J. Exp. Biol.* **204**, 1757–1764 (2001).
- Vellai, T., McCulloch, D., Gems, D. & Kovács, A. L. Effects of sex and insulin/insulin-like growth factor-1 signaling on performance in an associative learning paradigm in *Caenorhabditis elegans*. *Genetics* **174**, 309–316 (2006).
- Srinivasan, J. et al. A blend of small molecules regulates both mating and development in *Caenorhabditis elegans*. *Nature* **454**, 1115–1118 (2008).
- Barrios, A., Nurrish, S. & Emmons, S. W. Sensory regulation of *C. elegans* male mate-searching behavior. *Curr. Biol.* **18**, 1865–1871 (2008).
- Barr, M. M. & Sternberg, P. W. A polycystic kidney-disease gene homologue required for male mating behaviour in *C. elegans*. *Nature* **401**, 386–389 (1999).
- Barr, M. M. & García, L. R. Male mating behavior. *WormBook*. <http://dx.doi.org/10.1895/wormbook.1.78.1> (2006).

25. Lipton, J., Kleemann, G., Ghosh, R., Lints, R. & Emmons, S. W. Mate searching in *Caenorhabditis elegans*: a genetic model for sex drive in a simple invertebrate. *J. Neurosci.* **24**, 7427–7434 (2004).
26. Janssen, T. *et al.* Discovery and characterization of a conserved pigment dispersing factor-like neuropeptide pathway in *Caenorhabditis elegans*. *J. Neurochem.* **111**, 228–241 (2009).
27. Barrios, A., Ghosh, R., Fang, C., Emmons, S. W. & Barr, M. M. PDF-1 neuropeptide signaling modulates a neural circuit for mate-searching behavior in *C. elegans*. *Nature Neurosci.* **15**, 1675–1682 (2012).
28. Egger, B., Chell, J. M. & Brand, A. H. Insights into neural stem cell biology from flies. *Phil. Trans. R. Soc. Lond. B* **363**, 39–56 (2008).
29. Oikonomou, G. & Shaham, S. The glia of *Caenorhabditis elegans*. *Glia* **59**, 1253–1263 (2011).
30. Hong, Y., Roy, R. & Ambros, V. Developmental regulation of a cyclin-dependent kinase inhibitor controls postembryonic cell cycle progression in *Caenorhabditis elegans*. *Development* **125**, 3585–3597 (1998).
31. Zarkower, D. Somatic sex determination. *WormBook*. <http://dx.doi.org/10.1895/wormbook.1.84.1> (2006).
32. Hodgkin, J. A genetic analysis of the sex-determining gene, *tra-1*, in the nematode *Caenorhabditis elegans*. *Genes Dev.* **1**, 731–745 (1987).
33. Procko, C., Lu, Y. & Shaham, S. Sensory organ remodeling in *Caenorhabditis elegans* requires the zinc-finger protein ZTF-16. *Genetics* **190**, 1405–1415 (2012).
34. Hao, L., Johnsen, R., Lauter, G., Baillie, D. & Bürglin, T. R. Comprehensive analysis of gene expression patterns of hedgehog-related genes. *BMC Genomics* **7**, 280 (2006).
35. White, J. Q. *et al.* The sensory circuitry for sexual attraction in *C. elegans* males. *Curr. Biol.* **17**, 1847–1857 (2007).
36. Lee, K. & Portman, D. S. Neural sex modifies the function of a *C. elegans* sensory circuit. *Curr. Biol.* **17**, 1858–1863 (2007).
37. Mowrey, W. R., Bennett, J. R. & Portman, D. S. Distributed effects of biological sex define sex-typical motor behavior in *Caenorhabditis elegans*. *J. Neurosci.* **34**, 1579–1591 (2014).
38. WormBase. AMsoR http://www.wormbase.org/species/all/anatomy_term/WBbt:0003929#01-10.
39. Jarriault, S., Schwab, Y. & Greenwald, I. A *Caenorhabditis elegans* model for epithelial–neuronal transdifferentiation. *Proc. Natl Acad. Sci. USA* **105**, 3790–3795 (2008).
40. Okada, T. S. *Transdifferentiation: Flexibility in Cell Differentiation* (Clarendon Press, 1991).
41. Eguchi, G. & Kodama, R. Transdifferentiation. *Curr. Opin. Cell Biol.* **5**, 1023–1028 (1993).
42. Noctor, S. C., Martínez-Cerdeño, V., Ivic, L. & Kriegstein, A. R. Cortical neurons arise in symmetric and asymmetric division zones and migrate through specific phases. *Nature Neurosci.* **7**, 136–144 (2004).
43. Doetsch, F. The glial identity of neural stem cells. *Nature Neurosci.* **6**, 1127–1134 (2003).
44. Ninkovic, J. & Gotz, M. Fate specification in the adult brain—lessons for eliciting neurogenesis from glial cells. *Bioessays* **35**, 242–252 (2013).
45. Zurn, S. *et al.* Sequential histone-modifying activities determine the robustness of transdifferentiation. *Science* **345**, 826–829 (2014).
46. Arnold, A. P. Developmental plasticity in neural circuits controlling birdsong: sexual differentiation and the neural basis of learning. *J. Neurobiol.* **23**, 1506–1528 (1992).
47. WormAtlas (2002–2015) (eds Altun, Z. F. *et al.*) <http://www.wormatlas.org>.

Acknowledgements We would like to acknowledge M. Barr, in whose laboratory A.B. discovered the MCMs; WormAtlas for illustrations (reproduced with permission); T. Jarrell for contributions to the EM reconstruction; and W. Letton for the generation of strains and preliminary ablation studies. We thank M. Boxem, D. Portman, H. Baylis, L. Bianchi and R. Garcia for strains and reagents; M. Zhen, O. Hobert, I. Carrera, N. Stefanakis and S. Shaham, for unpublished reagents. Purified ascariosides were a gift from F. Schroeder to the Barr laboratory. Additional strains were obtained from the CGC, which is funded by NIH grant P40 OD010440. We thank L. Cochella, I. Carrera, S. Jarriault, and several of our close colleagues in CDB and NPP at University College London for discussions and comments on the manuscript; C. Barnes for advice on statistical analysis. This work was supported by a Master it! Scholarship Scheme (Malta and EU) to M.S., by NIH grant OD010943 to D.H.H., by Marie Curie CIG grant 618779 to R.J.P. and by a grant from The G. Harold and Leila Y. Mathers Charitable Foundation to S.W.E.; S.J.C. is supported by NIH grant 5T32GM007491; R.J.P. is a Wellcome Trust Research Career Development Fellow 095722/Z/11/Z; A.B. is supported by the Wellcome Trust Institutional Strategic Support Fund 097815/Z/11/A.

Author Contributions M.S., T.F., R.J.P. and A.B. conceived and performed the development and behaviour experiments. S.J.C., K.C.Q.N., S.W.E. and D.H.H. performed the ultrastructural analysis of the MCMs. S.J.C. and S.W.E. reconstructed the connectivity of the MCMs from serial EM sections. R.J.P. and A.B. co-wrote the manuscript and discussed it with all the authors.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to A.B. (a.barrios@ucl.ac.uk) or R.J.P. (r.poole@ucl.ac.uk).

METHODS

No statistical methods were used to predetermine sample size. No samples were excluded from analysis. For behavioural assays, investigators were blinded to the genotype or surgery of the animals during experiments and outcome assessment. Experiments testing response to secreted pheromones were randomized.

Strains. *him-5(e1490)* was used as the wild-type background strain. BL5717, *inIs179(ida-1::gfp);him-8(e1489)* was used for MCM ablation and behavioural assays (data in Figs 3 and 4 and Extended Data Fig. 2). BAR4, *otIs356(rab-3::rfp);him-5(e1490)* was used for MCM ablation and behavioural assays (data in Fig. 4 and Extended Data Fig. 2). PT2248, *pdf-1(tm1996);him-5(e1490)*. DR476, *daf-22(m130)*. CB369, *unc-51(e369)* VT774, *unc-36(e251);mIs103(rnr-1::gfp+unc-36(+))*. CB2823, *tra-1(e1488)III;eDp6(III);f*. BAR63, *stIs10116[his-72::his-24::mCherry::let-858'3'UTR+unc-119(+)]*; *itIs37[pie-1::mCherry::H2B::pie-1'3'UTR+unc-119(+)]*; *him-5(e1490)*.

Arrays for neuronal and amphid socket markers are indicated in Extended Data Tables 1 and 5, respectively.

DNA constructs. Sex transformation constructs were made by PCR fusion of the regulatory sequences of *grl-2* (ref. 34) or *ztf-16* (ref. 33) directly upstream of the ATG of *fem-3::SL2::mCherry* or *tra-2IC::SL2::mCherry*^{18,35–37}. These constructs include the *unc-54'3'UTR*. The *grl-2* promoter included 862 bp upstream of the ATG start codon (this is a smaller region of the promoter from that used in ref. 34). We checked transgenes driven by this promoter for embryonic expression and saw no expression before the birth of the AMso. The *ztf-16* glia enhancer included the region –2,536 to –4,637 upstream of the ATG, as described in ref. 33.

The following arrays were generated: *oleEx18[grl-2::fem-3::SL2mCherry(20 ng μ l⁻¹) + elt-2::gfp(40 ng μ l⁻¹)]*; *oleEx24[grl-2::fem-3::SL2mCherry(8 ng μ l⁻¹) + elt-2::gfp(40 ng μ l⁻¹)]*; *oleEx19[grl-2::tra-2IC::SL2mCherry(50 ng μ l⁻¹) + elt-2::gfp(30 ng μ l⁻¹)]*; *oleEx23[grl-2::tra-2IC::SL2mCherry(10 ng μ l⁻¹) + elt-2::gfp(40 ng μ l⁻¹)]*; *oleEx22[ztf-16::tra-2IC::SL2mCherry(5 ng μ l⁻¹) + elt-2::gfp(40 ng μ l⁻¹)]*.

Laser ablations. A standard protocol was used⁴⁸. L1, L3 or L4 males were mounted on a glass slide on a 5% agarose pad with 25 mM Na₂S₂O₃ as anaesthetic. For MCM candidate progenitor ablations the presence or absence of MCMs was identified with a *rab-3* transgene (Extended Data Table 1). Epithelial cells were identified by morphology, position, and expression of the *gals220[col-93::HIS-24::mCherry+unc-119(+)]* or *syIs78[ajm-1::gfp+unc-119(+)]* transgenes. Amphid socket cells were identified with glial markers (Extended Data Table 5).

For MCM ablation and behavioural assays the strains BL5717, *inIs179(ida-1::GFP);him-8(e1489)* and BAR4, *otIs356(rab-3::rfp);him-5(e1490)* were used. Animals were left to recover for one day and then assayed.

Behavioural assays. *Sexual conditioning switch.* Assays were performed as in ref. 10 with some modifications. Assay plates and 50 mM NaCl conditioning plates were 5 cm in diameter and contained 4.5 ml of 2% agar, 5 mM potassium phosphate (pH 6.0), 1 mM CaCl₂ and 1 mM MgSO₄ (salt-conditioning plates also contained 50 mM NaCl). A radial gradient of salt was created in the assay plates by adding 10 μ l of 50 mM NaCl 2 cm away from the centre of the plate the night before the assay and then 5 μ l of 50 mM NaCl three hours before the assay. Males and hermaphrodites were picked as L4s the night before the assay and transferred to single-sex plates with food. Animals were recovered from the food plates with wash buffer (1 mM CaCl₂, 1 mM MgSO₄, and 5 mM pH 6.0 potassium phosphate) and centrifuged at 1,700 r.p.m. for 3 min. A total of three washes were performed to remove the food before placing the animals on the conditioning plates. Between 8 and 15 males were placed on each conditioning plate and 200 hermaphrodites were used for sexual conditioning. Animals were conditioned for 2 h. Ablated males were sexually conditioned in the same plate as control males, tested blind, and identified afterwards based on the absence of MCMs. Animals were tested individually on a radial gradient of salt by placing them in the centre of the assay plate, 2 cm away from the source of NaCl. After 20 min, the tracks left by the animal were visualized and the chemotaxis score was calculated as in ref. 10 by the sum of the scores of the regions by which the animal had travelled (Fig. 3a). At least five replicas were performed and on different days. Mann–Whitney *U*-test was used for statistical analysis.

Response to secreted mate pheromones. Assays were performed as in ref. 22. Hermaphrodite-conditioned media was prepared by incubating wild-type or *daf-22(m130)* hermaphrodites in M9 buffer at a concentration of one hermaphrodite per 1 μ l of M9 for 3 h. Animals were picked as L4 s the night before. Assay plates were regular NGM plates seeded with 50 μ l of OP50. A pheromone spot (1 μ l of hermaphrodite-conditioned M9, or dilutions of Asc no. 2 and Asc no. 3 in M9) and a control spot (1 μ l of M9) were placed 3 mm apart. The two spots were interchanged every other trial to remove any bias. Upon drying, the spots left a visible rim on the food lawn. A population of 5 males (8 males for *him-5(e1490)* and *pdf-1(tm1996)* trials) was placed 3 mm away from the spots and video

recorded for 20 min. Videos were scored blindly for time spent in each spot. Each entrance to a spot was considered an independent event if after leaving, the worm had moved at least 2 mm away from the spot. The average duration of events was calculated and compared between spots and populations. Mann–Whitney *U*-test with Bonferroni correction was used for statistical analysis. For clear visualization of the data in the plotted graphs, data was logarithmically transformed. This results in normalization of the distribution of the data. *t*-test with Bonferroni correction was used for statistical analysis of the transformed data. Comparison of arithmetic (raw data) or geometric (transformed data) means did not affect the interpretation of the data.

Seven independent population trials were performed for each condition except for *daf-22(m130)*-conditioned media (twelve), 600 nM Asc no. 2 + 60 nM Asc no. 3 (five) and 1200 nM Asc no. 2 + 120 nM Asc no. 3 (four). We selected three different ascaroside dilutions according to the response elicited in intact males: 600 nM Asc no. 2 + 60 nM Asc no. 3 (no robust and stochastic response), 800 nM Asc no. 2 + 80 nM Asc no. 3 (robust but low response), 1,200 nM Asc no. 2 + 120 nM Asc no. 3 (robust and strong response) (Extended Data Fig. 2).

Response efficiency to mate contact. Animals were isolated as L4s the night before. For the assay, single males were placed with 30 mates (hermaphrodites) in a 20 mm food lawn. Males were tested until they responded to a mate or for three minutes, whatever happened first. A male was scored as responding to mate contact if it placed its tail ventral down on the mate's body and initiated the mating sequence by backing along the mate's body to make a turn. The response efficiency was calculated by dividing a response by the total number of contacts made with the mate before responding. If a male did not respond within three minutes, it was scored as having 0 efficiency. Males were scored blindly. The Mann–Whitney *U*-test was used for statistical analysis.

Turning and location of vulva (lov) assays. Single males (isolated the night before as L4) were tested on a 10 μ l lawn of food (*E. coli* OP50) with 30 *unc-51* hermaphrodites during 3 min. Upon response to contact, up to seven turns were scored as good or bad. Good turns: sharp ventral bend around the tip of the body of the mate, followed by uninterrupted scanning along the other side of the mate's body. Bad turns: stutter (stops moving backwards before turning and moves forward before continuing with backward move); wide (tail positioned in a loose bend when turning around the tip of the mate); swim off (continues backing after reaching the end of the mate's body and loses contact); or interrupted (after turning, the male does not continue backing along the mate's body). The proportion of good turns out of total turns was scored.

For location of vulva efficiency, the number of encounters with the vulva until it was first located was scored. The assay finished when the male located the vulva or after three minutes, whatever happened first. Location efficiency was scored as 1 divided by the total number of encounters.

The Mann–Whitney *U*-test was used for statistical analysis.

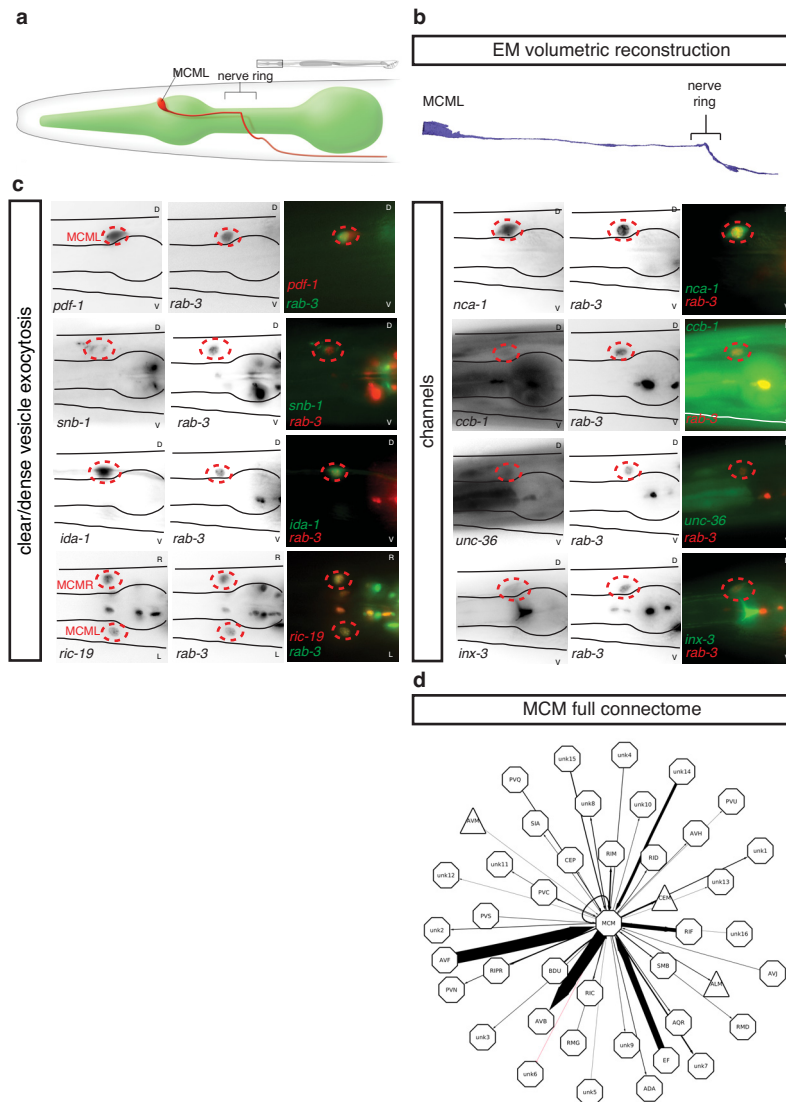
Fertility assays. A single *rab-3::rfp* transgenic male was left to mate with two *unc-51* hermaphrodites on a 15- μ l lawn of food (*E. coli* OP50) for 24 h and then removed. The hermaphrodites were transferred to a fresh plate every 24 h for a total of three days. The proportion of cross-progeny from total progeny was calculated for three ablated individuals and three controls. *t*-test was used for statistical analysis.

Food-leaving and retention assays. As described in ref. 27.

Electron microscopy and serial reconstruction. Four samples were fixed by chemical fixation or high-pressure freezing and freeze substitution as previously described⁴⁹. Ultrathin sections were cut using a RMC PowerTome XL, collected onto grids, and imaged using either a Philips CM10 TEM or Zeiss Supra 40 FE-SEM. Sections were elastically aligned⁵⁰ and volumetrically reconstructed using TrakEM2 (ref. 51). The MCM cell bodies were identified in the EM sections based on position and morphology and in comparison to similar hermaphrodite sections. This identity was further confirmed by the identification of a posterior projection, as no other cell with its soma in the same region is known to project posteriorly. This was followed by serial tracing of the projections to establish their morphology and connectivity. Synaptic connectivity and skeleton diagrams were determined using Elegance⁵². Circuit diagrams of connectivity were generated using Cytoscape⁵³. We estimated the anatomical strength of synaptic connectivity between two neurons by summing the number of serial sections where we observed the ultrastructural morphology presynaptic components using the same criteria as ref. 13.

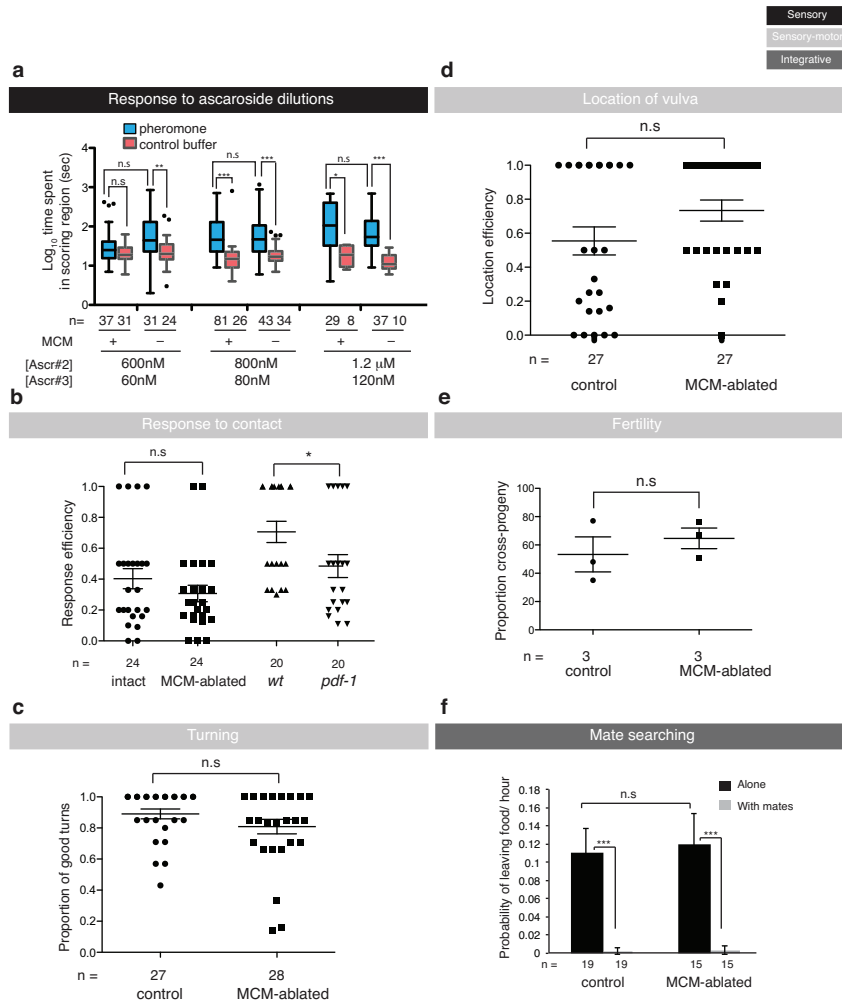
48. Bargmann, C. I. & Avery, L. Laser killing of cells in *Caenorhabditis elegans*. *Methods Cell Biol.* **48**, 225–250 (1995).
49. Hall, D. H., Hartwig, E. & Nguyen, K. C. Q. Modern electron microscopy methods for *C. elegans*. *Methods Cell Biol.* **107**, 93–149 (2012).
50. Saalfeld, S., Fetter, R., Cardona, A. & Tomancak, P. Elastic volume reconstruction from series of ultra-thin microscopy sections. *Nature Methods* **9**, 717–720 (2012).

51. Cardona, A. *et al.* TrakEM2 software for neural circuit reconstruction. *PLoS ONE* **7**, e38011 (2012).
52. Xu, M. *et al.* Computer assisted assembly of connectomes from electron micrographs: application to *Caenorhabditis elegans*. *PLoS ONE* **8**, e54050 (2013).
53. Smoot, M. E., Ono, K., Ruscheinski, J., Wang, P.-L. & Ideker, T. Cytoscape 2.8: new features for data integration and network visualization. *Bioinformatics* **27**, 431–432 (2011).
54. Tursun, B., Patel, T., Kratsios, P. & Hobert, O. Direct conversion of *C. elegans* germ cells into specific neuron types. *Science* **331**, 304–308 (2011).
55. Altun-Gultekin, Z. *et al.* A regulatory cascade of three homeobox genes, *ceh-10*, *ttx-3* and *ceh-23*, controls cell fate specification of a defined interneuron class in *C. elegans*. *Development* **128**, 1951–1969 (2001).
56. Barrios, A., Ghosh, R., Fang, C., Emmons, S. W. & Barr, M. M. PDF-1 neuropeptide signaling modulates a neural circuit for mate-searching behavior in *C. elegans*. *Nature Neuroscience* **15**, 1675–1682 (2012).
57. Zahn, T. R., Macmorris, M. A., Dong, W., Day, R. & Hutton, J. C. IDA-1, a *Caenorhabditis elegans* homolog of the diabetic autoantigens IA-2 and phogrin, is expressed in peptidergic neurons in the worm. *J. Comp. Neurol.* **429**, 127–143 (2000).
58. Stefanakis, N., Carrera, I. & Hobert, O. Regulatory logic of pan-neuronal gene expression in *C. elegans*. *Neuron* **87**, 733–750 (2015).
59. McKay, S. J. *et al.* Gene expression profiling of cells, tissues, and developmental stages of the nematode *C. elegans*. *Cold Spring Harb. Quant. Biol.* **68**, 159–169 (2015).
60. Altun, Z. F., Chen, B., Wang, Z. W. & Hall, D. H. High resolution map of *Caenorhabditis elegans* gap junction proteins. *Developmental Dynamics* **238**, 1936–1950 (2009).
61. Serrano-Saiz, E. *et al.* Modular control of glutamatergic neuronal identity in *C. elegans* by distinct homeodomain proteins. *Cell* **155**, 659–673 (2013).
62. Kratsios, P., Stolfi, A., Levine, M. & Hobert, O. Coordinated regulation of cholinergic motor neuron traits through a conserved terminal selector gene. *Nature Neuroscience* **15**, 205–214 (2011).
63. Bell, L. R., Stone, S., Yochem, J., Shaw, J. E. & Herman, R. K. The molecular identities of the *Caenorhabditis elegans* intraflagellar transport genes *dyl-6*, *daf-10* and *osm-1*. *Genetics* **173**, 1275–1286 (2006).
64. Gerisch, B., Weitzel, C., Kober-Eisermann, C., Rottiers, V. & Antebi, A. A hormonal signaling pathway influencing *C. elegans* metabolism, reproductive development, and life span. *Dev. Cell* **1**, 841–851 (2001).
65. Alkema, M. J., Hunter-Ensor, M., Ringstad, N. & Horvitz, H. R. Tyramine functions independently of octopamine in the *Caenorhabditis elegans* nervous system. *Neuron* **46**, 247–260 (2005).
66. Kim, K. & Li, C. Expression and regulation of an FMRFamide-related neuropeptide gene family in *Caenorhabditis elegans*. *J. Comp. Neurol.* **475**, 540–550 (2004).
67. Beets, I. *et al.* Vasopressin/oxytocin-related signaling regulates gustatory associative learning in *C. elegans*. *Science* **338**, 543–545 (2012).
68. Gruninger, T. R., Gualberto, D. G., LeBoeuf, B. & Garcia, L. R. Integration of male mating and feeding behaviors in *Caenorhabditis elegans*. *J. Neurosci.* **26**, 169–179 (2006).
69. Yoshimura, S., Murray, J. I., Lu, Y., Waterston, R. H. & Shaham, S. *mls-2* and *vab-3* control glia development, *hlh-17*/Olig expression and glia-dependent neurite extension in *C. elegans*. *Development* **135**, 2263–2275 (2008).
70. Haklai-Topper, L. *et al.* The neurexin superfamily of *Caenorhabditis elegans*. *Gene Expr. Patterns* **11**, 144–150 (2011).
71. Gower, N. J. D. *et al.* Dissection of the promoter region of the inositol 1,4,5-trisphosphate receptor gene, *itr-1*, in *C. elegans*: a molecular basis for cell-specific expression of IP₃R isoforms. *J. Mol. Biol.* **306**, 145–157 (2011).
72. Heiman, M. G. & Shaham, S. DEX-1 and DYF-7 establish sensory dendrite length by anchoring dendritic tips during cell migration. *Cell* **137**, 344–355 (2009).



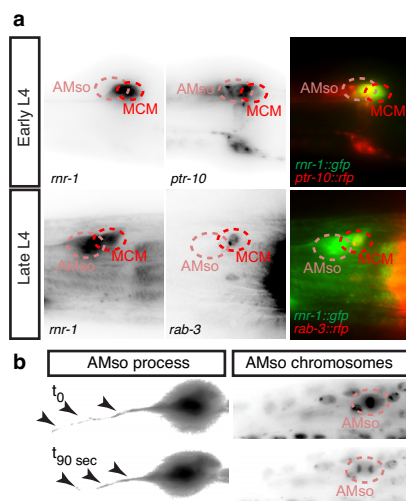
Extended Data Figure 1 | The MCMs are newly identified male-specific neurons. **a**, WormAtlas-style diagram depicting the morphology and position of one of the bilateral pair of MCM neurons in the head of a male worm and its projection within the nerve ring and along the ventral cord. **b**, Volumetric reconstruction of the MCML cell body and projection based on tracing of serial EM sections. **c**, Co-expression of transgenes for neuronal markers in the *rab-3*-positive cells identified as MCMs (indicated with dashed red circles). All photographs are lateral views of animals oriented anterior to the left and dorsal to the top except for *ric-19*, which are dorsal views. Transgenes are listed in

Extended Data Table 1. *pdf-1* (neuropeptide pigment dispersing factor); *snb-1* (synaptobrevin); *ida-1* (tyrosine phosphatase-like receptor, orthologue of mammalian phogrin); *ric-19* (rab-2 effector); *nca-1* (NALCN Na⁺ channel subunit); *ccb-1* (voltage-gated Ca²⁺ channel subunit); *unc-36* (voltage-gated Ca²⁺ channel subunit); *inx-3* (gap junction innexin). D, dorsal; L, left; R, right; V, ventral. **d**, Diagram of the neurons that directly connect to and from the MCMs. Triangles, sensory neurons; octagons, interneurons and unidentical neurons. The thickness of the arrows is proportional to the anatomical strength of the connections (Extended Data Table 2).

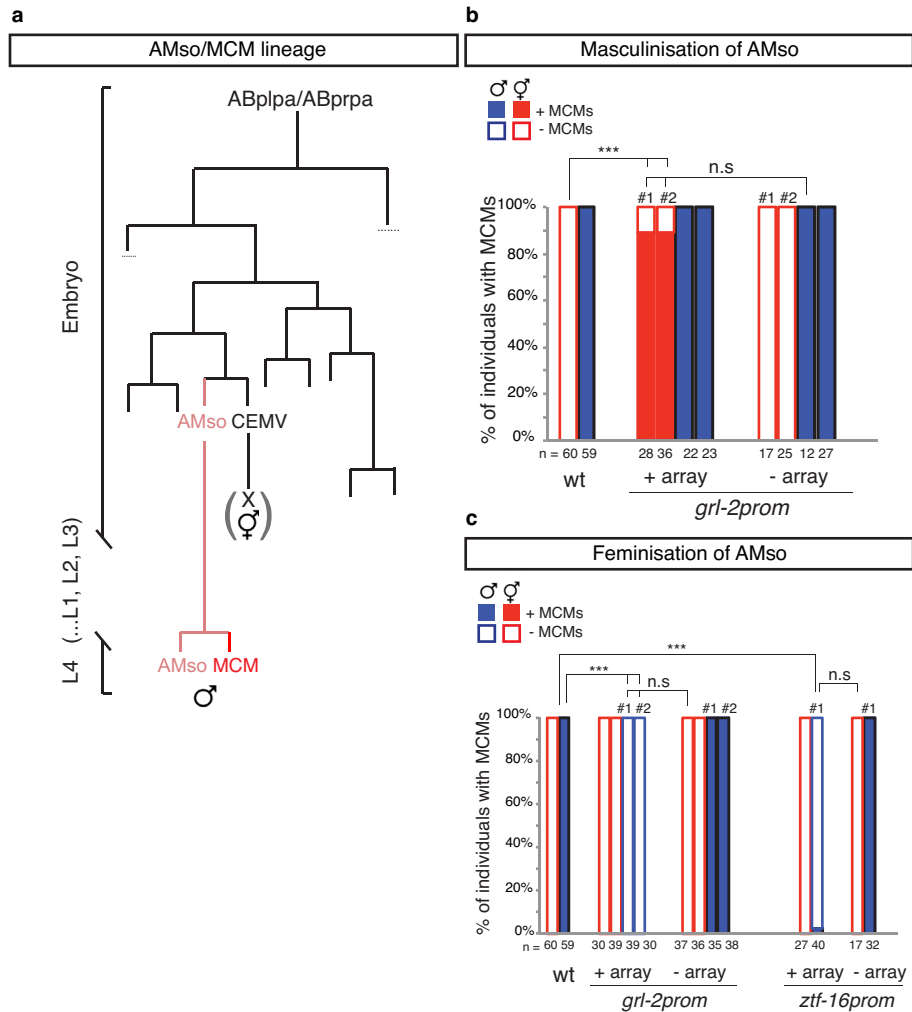


Extended Data Figure 2 | The MCMs are not required for other male-specific behaviours. **a**, Response of intact and MCM-ablated males (*inIs179(ida-1::gfp);him-8(e1489)* and *otIs356(rab-3::rfp)him-5(e1490)*) to dilutions of ascaroside pheromones (Ascr). Graphs represent Tukey box plots of logarithmic transformations of the data; *n*, number of independent events (that is, entry in scoring region). *t*-test with Bonferroni correction was used for statistical analysis. ****P* < 0.001; ***P* < 0.01; **P* < 0.05; n.s., no statistically significant difference (*P* ≥ 0.05). **b**, Response efficiency to mate contact of intact, MCM-ablated and *pdf-1(tm1996)* mutant males measured as the proportion of responses out of total contacts with an hermaphrodite. Intact and MCM-ablated animals were *inIs179(ida-1::gfp);him-8(e1489)*. Wild-type animals were *him-5(e1490)*. A response indicates that the male placed its tail ventral down on the mate's body and backed along it to make a turn.

c, d, Proportion of good turns (**c**) and location of vulva efficiency (**d**) of intact and MCM-ablated males (*inIs179(ida-1::gfp);him-8(e1489)* and *otIs356(rab-3::rfp)him-5(e1490)*). **e**, Fertility (measured as proportion of cross-progeny) of intact and MCM-ablated males (*otIs356(rab-3::rfp)him-5(e1490)*). For **b–e**, *n*, number of individual animals tested. Error bars indicate s.e.m. Mann–Whitney *U*-test was used for statistical analysis. **P* < 0.05; n.s., no statistically significant difference (*P* ≥ 0.05). **f**, Mate-searching behaviour, measured as *P_L* values (probability of leaving food per hour) in the absence or presence of mates, of intact and MCM-ablated males (*otIs356(rab-3::rfp)him-5(e1490)*). *n*, number of individual animals tested. Two independent population assays were performed on different days. Maximum likelihood statistical analysis was used to compare *P_L* values. Error bars indicate s.e.m. ****P* < 0.001; n.s., no statistically significant difference (*P* ≥ 0.05).



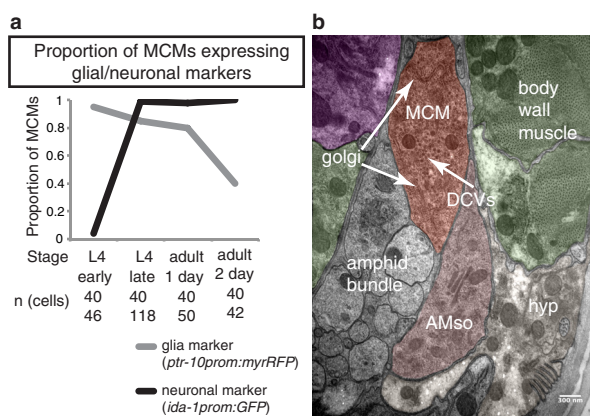
Extended Data Figure 3 | The MCMs arise from a division of the AMso glial cell. All photographs are lateral views of animals oriented anterior to the left and dorsal to the top. **a**, Fluorescent photographs showing the two cells expressing *mr-1::gfp* co-labelled with the glial marker *ptr-10::rfp* and the neuronal marker *rab-3::rfp* in the head of males at the early and late L4 stages. The AMso and MCM cell bodies are indicated with dashed lines. **b**, Fluorescent images of the AMso cell body and its projection at two time points during cell division. Photos are overexposed for visualization of the projection, indicated by arrows. The chromosomes are labelled with a histone::*rfp* transgene, and the AMso cell body is indicated by dashed lines.



Extended Data Figure 4 | AMso plasticity is regulated by AMso genetic sex.

a, Diagram of the AMso and MCM lineage. **b**, **c**, Proportion of individuals with MCMs in control animals and animals expressing sex-reversing transgenes in AMso. **b**, AMso masculinization with *grl-2::fem-3::SL2::mCherry* transgenes (*oleEx18* and *oleEx24*). **c**, AMso feminization with *grl-2::tra-2IC::SL2::mCherry* transgenes (*oleEx19* and *oleEx23*) and *ztf-16::tra-2IC::SL2::mCherry* transgene *oleEx22*. MCM cell fate was identified with *ida-1::gfp* or

rab-3::yfp reporter transgenes. In the head, the *grl-2* promoter drives expression in AMso and the excretory duct and pore cells, and the *ztf-16* glial enhancer drives expression in the AMso and amphid sheath glia. # indicates an independent transgenic array line for each manipulation. χ^2 test was used for statistical analysis; *** $P < 0.001$; n.s., no statistical significant difference ($P \geq 0.05$); n = number of animals scored.



Extended Data Figure 5 | The MCMs lose molecular and structural

characteristics of glia after birth. **a**, Proportion of MCMs with presence of the glial marker *ptr-10::myrRfp* or the neuronal marker *ida-1::gfp* at different stages after MCM birth. **b**, Electron micrograph of a cross-section of an adult male head showing the MCM and AMso cell body ultrastructure. Neighbouring tissues are colour coded following WormAtlas (<http://www.wormatlas.org/colorcode.htm>). Purple (pharynx), muscle (green), hypodermis (light cream), AMso (amphid socket, pink). The dendrites of the amphid neurons (amphid bundle) are not colored.

Extended Data Table 1 | Reporter transgenes for neuronal markers tested for MCM expression

protein/function	gene	array	MCM expression
Neuronal GTPase	<i>rab-3</i>	<i>otIs291 (rab-3prom:: NLS YFP+rol-6) and otIs356(rab-3prom::rfp)</i> ⁵⁴	+ 62/62
Neuronal GEF	<i>rgef-1</i>	<i>evIs111 (rGEFprom:: GFP)</i> ⁵⁵	+
Neuropeptide	<i>pdf-1</i>	<i>myEx696 (pdf-1prom::RFP+unc-122::GFP)</i> ⁵⁶	+
Phogrin/Exocytosis	<i>ida-1</i>	<i>inIs179(ida-1prom::GFP)</i> ⁵⁷	+ 59/59
Neuropeptide exocytosis	<i>ric-19</i>	<i>otEx6173 (ric-19prom6::NLS-TagRFP)</i> (gift from Hobert lab) ⁵⁸	+
Syntaxin	<i>unc-64</i>	<i>otEx4553 (fosmid-based transcriptional reporter)</i> (gift from Hobert lab) ⁵⁸	+
Synaptobrevin	<i>snb-1</i>	<i>otEx5422 (fosmid-based translational reporter)</i> (gift from Hobert lab)	+
Synaptotagmin	<i>snt-1</i>	<i>otEx5924(fosmid-based transcriptional reporter)</i> (gift from Hobert lab) ⁵⁸	+
Na ²⁺ channel subunit	<i>nca-1</i>	<i>hpEx528 (nca-1prom::GFP)</i> (gift from Zhen lab)	+
Ca ²⁺ channel subunit	<i>ccb-1</i>	<i>sEX10836(T28F2.5prom::GFP+pCeh361)</i> ⁵⁹	+
Ca ²⁺ channel subunit	<i>unc-36</i>	<i>sEX12299(C50C3.9prom::GFP+pCeh361)</i> ⁵⁹	+
Innexin	<i>inx-3</i>	<i>zwEx103 [inx-3prom::GFP +lin-15(+)]</i> ⁶⁰	+
VGLUT	<i>eat-4</i>	<i>otIs388 (fosmid-based reporter)</i> ⁶¹	–
Choline transporter	<i>cho-1</i>	<i>otIs323 (fosmid-based reporter)</i> ⁶²	–
VGAT	<i>unc-47</i>	<i>otIs348[unc-47prom(300bp)::mCHOPTI::unc54-3'UTR), pha-1(+)]</i> ⁶²	–
Intraflagellar transport (marker for ciliated neurons)	<i>osm-6</i>	<i>mnIs17(osm-6prom::osm-6:GFP)</i> ⁶³	–
Fatty acid hydroxylase (marker for the XXX cell)	<i>daf-9</i>	<i>dhIs59(daf-9prom::GFP)</i> ⁶⁴	–
Octopamine synthesis	<i>tbh-1</i>	<i>nIs107[tbh-1prom::GFP+lin-15(+)]</i> ⁶⁵	–
Neuropeptide	<i>flp-21</i>	<i>nyIs80 (flp-21prom::GFP)</i> ⁶⁶	–
Nematocin receptor	<i>ntr-1</i>	<i>lStEx33 [ntr-1prom::gfp]</i> ⁶⁷	–
Ca ²⁺ channel subunit	<i>cca-1</i>	<i>sEX14060(C54D2.5aprom::GFP+pCeh361)</i> ⁵⁹	–
K ⁺ channel subunit	<i>unc-103</i>	<i>rgEx[unc-103prom::YFP + pBx1]</i> ⁶⁸	–
Innexin	<i>inx-12</i>	<i>zwEx112 [inx-12prom::GFP +lin-15(+)]</i> ⁶⁰	–
Innexin	<i>inx-13</i>	<i>zwEx113 [inx-13prom::GFP +lin-15(+)]</i> ⁶⁰	–
Innexin	<i>inx-5</i>	<i>zwEx105 [inx-5prom::GFP +lin-15(+)]</i> ⁶⁰	–
Innexin	<i>inx-12</i>	<i>zwEx112 [inx-12prom::GFP +lin-15(+)]</i> ⁶⁰	–

+ indicates presence; – indicates absence. For integrated arrays, quantification is indicated. At least 15 animals were examined for each transgene. This Table cites refs 54–68.

Extended Data Table 2 | MCM connectivity

type	presynaptic	postsynaptic	#synapses	#EM sections
chemical	AVF	MCM	16	61
chemical	AVJ	MCM	1	2
chemical	AVM	MCM	1	1
chemical	BDU	MCM	3	7
chemical	CEM	MCM	2	5
chemical	CEP	MCM	1	1
chemical	EF	MCM	6	35
chemical	MCM	ADA	1	3
chemical	MCM	ALM	2	3
chemical	MCM	AQR	2	3
chemical	MCM	AVB	29	61
chemical	MCM	AVF	1	2
chemical	MCM	AVH	1	2
chemical	MCM	BDU	1	3
electrical	MCM	CEM	1	1
chemical	MCM	CEM	5	8
chemical	MCM	CEP	1	1
chemical	MCM	EF	2	3
chemical	MCM	MCM	2	6
chemical	MCM	PVC	1	2
chemical	MCM	PVN	3	3
chemical	MCM	PVU	1	1
chemical	MCM	RIC	2	3
chemical	MCM	RID	1	3
chemical	MCM	RIF	11	23
chemical	MCM	RIM	3	7
chemical	MCM	RIPR	3	6
chemical	MCM	RMD	1	2
chemical	MCM	SIA	1	2
chemical	MCM	SMB	2	3

type	presynaptic	postsynaptic	#synapses	#EM sections
chemical	MCM	unk1 ⁺	2	4
chemical	MCM	unk10 ⁺	1	2
chemical	MCM	unk11 ⁺	1	2
chemical	MCM	unk12 ⁺	1	1
chemical	MCM	unk13 ⁺	1	1
chemical	MCM	unk2 ⁺	1	3
chemical	MCM	unk3 ⁺	1	2
electrical	MCM	unk6 ⁺	1	1
chemical	MCM	unk7 ⁺	3	7
chemical	MCM	unk8 ⁺	4	4
chemical	MCM	unk9 ⁺	1	2
chemical	PVC	MCM	1	3
chemical	PVQ	MCM	1	4
chemical	PVS	MCM	1	2
chemical	RIC	MCM	1	1
chemical	RID	MCM	1	1
chemical	RIF	MCM	10	20
chemical	RIM	MCM	3	3
chemical	RMG	MCM	1	2
chemical	SMB	MCM	2	3
chemical	unk14 ⁺	MCM	3	15
chemical	unk15 ⁺	MCM	1	5
chemical	unk16 ⁺	MCM	1	1
chemical	unk4 ⁺	MCM	2	3
chemical	unk5 ⁺	MCM	1	1

*unk refers to neurons whose identity has not been confirmed unambiguously in the EM serial sections.

Extended Data Table 3 | Cell ablations of candidate MCM progenitors

ablated cell	stage	animals with loss of MCM/total
AMso	L3	6/7
AMso	L4	0/3
H0 seam cell	L1	1/16
H1 seam cell	L1	0/2
H0+H1	L1	0/1
Hyp1 (lateral nucleus)	L1	0/4
H0+Hyp1(lateral nucleus)	L1	0/1

Extended Data Table 4 | Mosaic analysis of sex-transformation arrays, scoring the presence of MCMs

sex reversal	array	array in AMso and other cells [*]	array in other cells [*] only
Masculinised hermaphrodites	<i>oleEx18 (grl-2prom::fem-3:SL2:mCherry)</i>	89% (n=28)	0% (n=22)
Feminised males	<i>oleEx19 (grl-2prom::tra-2IC:SL2:mCherry)</i>	0% (n=39)	100% (n=7)
Feminised males	<i>oleEx22 (ztf-16prom::tra-2IC:SL2:mCherry)</i>	2.5% (n=40)	87.5%(n=16)

^{*}Expression of *oleEx18 (grl-2::fem-3::SL2::mCherry)* in the head was observed in AMso, excretory duct and pore cells and sometimes in the pharynx and/or hypodermis. Expression of *oleEx19 (grl-2::tra-2IC::SL2::mCherry)* in the head was observed in AMso, excretory duct and pore cells and sometimes in the hypodermis. Expression of *oleEx22 (ztf-16::tra-2IC::SL2::mCherry)* in the head was observed in AMso and sometimes in the amphid sheath and/or a neuron in the nerve ring.

Extended Data Table 5 | Reporter transgenes for glial/AMso markers

protein/function	gene	array	glial subtype expression	AMso expression
Patched (PTCHD3) related receptor	<i>ptr-10</i> ⁶⁹	<i>nsIs108[ptr-10prom::myristyl-Rfp]</i>	All glia	+
Hedgehog-like/Ground-related	<i>grl-2</i> ^{3b}	<i>sEx12852[T16G1.8prom::GFP+pCeh361]</i>	AMso and PHso	+
Inositol trisphosphate receptor	<i>itr-1</i> ^{71, 72}	<i>jwEx51(itr1promB::GFP+rol-6) and nsEX1153 [F16F9.3prom::mCherry+itr-1prom::CFP+rol-6(su1006)]</i>	AMso	+
C2H2 zinc-finger transcription factor	<i>ztf-16</i> ³⁴	<i>oleEx22[ztf-16enhancer::tra-2IC::SL2cherry+elt-2::GFP]</i>	AMso, AMsh, PHso	+
Basic helix-loop-helix transcription factor	<i>hlh-17</i> ⁶⁹	<i>leEx1713[hlh-17prom::GFP+unc-119(+)]</i>	CEPsh	–
Caspr (Neurexin superfamily)	<i>itx-1</i> ⁷⁰	<i>otEx[W03D8.6prom::GFP+rol-6]</i>	OL and IL glia	–

+ indicates presence; – indicates absence. At least 15 animals were examined for each transgene. AMso, amphid socket; PHso, phasmid socket; AMsh, amphid sheath; PHsh, phasmid sheath; CEPsh, cephalic sensilla sheath; OLL, outer labial; IL, inner labial. This Table cites refs 34, 35, 69–72

Molecular basis of ligand recognition and transport by glucose transporters

Dong Deng^{1,2,3*}, Pengcheng Sun^{1,2,3*}, Chuangye Yan^{1,2,3}, Meng Ke^{1,2,3}, Xin Jiang^{1,2}, Lei Xiong³, Wenlin Ren^{1,2}, Kunio Hirata^{4,5}, Masaki Yamamoto⁴, Shilong Fan² & Nieng Yan^{1,2,3}

The major facilitator superfamily glucose transporters, exemplified by human GLUT1–4, have been central to the study of solute transport. Using lipidic cubic phase crystallization and microfocus X-ray diffraction, we determined the structure of human GLUT3 in complex with D-glucose at 1.5 Å resolution in an outward-occluded conformation. The high-resolution structure allows discrimination of both α - and β -anomers of D-glucose. Two additional structures of GLUT3 bound to the exofacial inhibitor maltose were obtained at 2.6 Å in the outward-open and 2.4 Å in the outward-occluded states. In all three structures, the ligands are predominantly coordinated by polar residues from the carboxy terminal domain. Conformational transition from outward-open to outward-occluded entails a prominent local rearrangement of the extracellular part of transmembrane segment TM7. Comparison of the outward-facing GLUT3 structures with the inward-open GLUT1 provides insights into the alternating access cycle for GLUTs, whereby the C-terminal domain provides the primary substrate-binding site and the amino-terminal domain undergoes rigid-body rotation with respect to the C-terminal domain. Our studies provide an important framework for the mechanistic and kinetic understanding of GLUTs and shed light on structure-guided ligand design.

Cellular uptake of glucose is a fundamental process for metabolism, growth and homeostasis¹. The SLC2 family glucose transporters (GLUTs) catalyse facilitative diffusion of glucose and other monosaccharides across biomembranes^{2,3}. The fourteen human GLUTs have specific spatial and temporal distributions and exhibit distinct transport kinetics, capacity and substrate selectivity⁴.

GLUT1–4 are among the most rigorously characterized solute transporters. GLUT1, which was one of the first transporters to be characterized^{5–7}, has been a paradigm in the understanding of solute transport^{1,8}. GLUT1 is the principal glucose transporter in erythrocytes and blood-tissue barriers, and ubiquitous for basal-level glucose uptake^{9,10}. GLUT2 represents the major hepatocyte isoform that allows both uptake and efflux of glucose in response to fed or fasted state. GLUT2 also mediates glucose transport in intestinal, kidney and β -pancreatic cells^{11,12}. GLUT3 is referred to as the ‘neuronal glucose transporter’ for its primary function in neurons, and it is also responsible for glucose uptake in sperm, preimplantation embryos and circulating white blood cells¹³. GLUT4 is responsive to insulin in adipocytes and muscles^{14,15}.

Inactivating mutations or mis-regulations of GLUTs are associated with deleterious diseases including GLUT1 deficiency syndrome (the De Vivo disease), Fanconi–Bickel syndrome, type 2 diabetes mellitus and Alzheimer’s disease^{16–19}. GLUT1 and GLUT3 are overexpressed in different types of solid tumours, where the demand for glucose is strikingly enhanced to compensate for ATP generation under anaerobic conditions (the Warburg effect)^{20–24}. Targeting the overexpression of GLUTs for cancer diagnosis and potential therapy has drawn increasing attention²⁵, exemplified by the positron emission tomography that monitors the uptake of 2-deoxy-2-[¹⁸F]fluoroglucose^{1,26,27}. Glucose conjugations were designed to use glucose transporters for enhanced membrane permeation and tissue-specific delivery of anti-cancer

drugs²⁸. Structural determination of GLUTs, particularly in complex with ligands, is a prerequisite for ligand design and optimization.

GLUTs function by the ‘alternating access’ mechanism, whereby the alternate exposure of the substrate-binding site(s) to either side of the membrane is achieved through cycles of conformational changes of the transporter^{29–31}. The structure of GLUT1 was recently captured in an inward-open conformation³². Structures of Xyle, the *Escherichia coli* homologue of GLUTs, are available in three states: ligand-bound and outward-occluded³³, inward-open^{34,35}, and partly inward-occluded³⁴.

Here we present the *in meso* crystal structure of human GLUT3 bound to D-glucose in an outward-occluded state at 1.5 Å resolution. Structures of GLUT3 in complex with the exofacial competitive inhibitor maltose were also obtained in the outward-open and outward-occluded conformations. Structural comparison with the inward-open GLUT1 reveals the molecular basis for substrate recognition and transport by GLUTs.

Characterizations of the recombinant GLUT3(N43T)

All the functional and structural characterizations of GLUT3 reported here were performed using the glycosylation-site-eliminated variant GLUT3(N43T), which will be referred to as GLUT3 hereafter. A proteoliposome-based counterflow assay was reconstituted following a modified protocol^{36,37}, and the substrate selectivity of the recombinant GLUT3 was examined qualitatively (Fig. 1a, b). Among fifteen tested hexoses and pentoses, D-glucose, D-galactose, D-mannose, D-xylose and D-fucose blocked more than 90% of the [³H]glucose uptake, and L-arabinose and D-lyxose inhibited more than half of the activity. The other sugars had little or no effect on the transport. Notably, the C3 hydroxyl (C3-OH) of the seven effective inhibitors exhibits the same equatorial configuration, suggesting that C3-OH

¹State Key Laboratory of Membrane Biology, Tsinghua University, Beijing 100084, China. ²Center for Structural Biology, Tsinghua University, Beijing 100084, China. ³Tsinghua-Peking Center for Life Sciences, School of Life Sciences and School of Medicine, Tsinghua University, Beijing 100084, China. ⁴Advanced Photon Technology Division, Research Infrastructure Group, SR Life Science Instrumentation Unit, RIKEN/SPring-8 Center, 1-1-1 Kouto Sayo-cho Sayo-gun, Hyogo 679-5148 Japan. ⁵Precursory Research for Embryonic Science and Technology (PRESTO), Japan Science and Technology Agency, 4-1-8 Honcho, Kawaguchi, Saitama 332-0012, Japan.

*These authors contributed equally to this work.

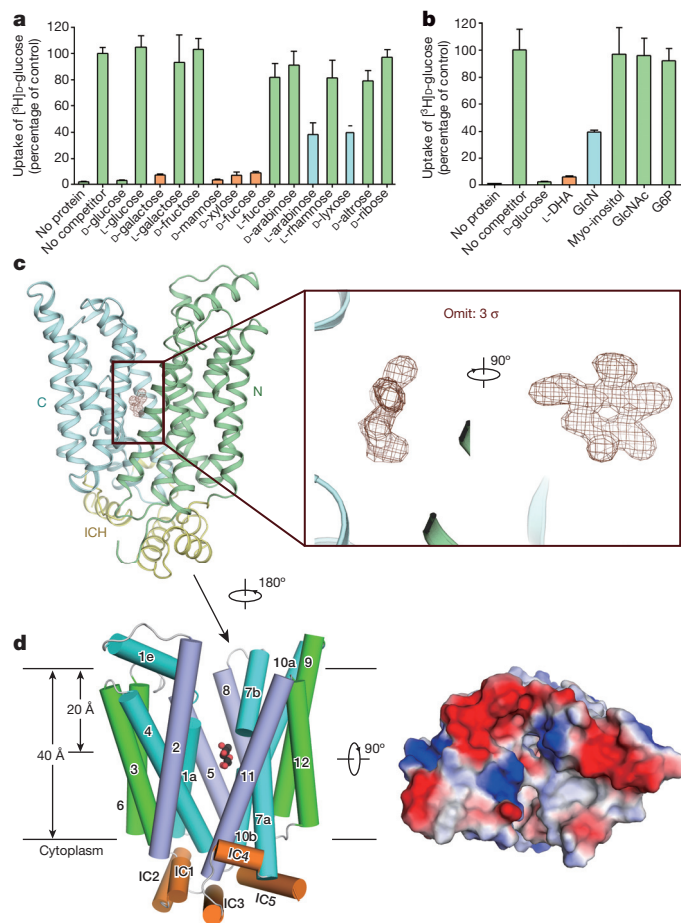


Figure 1 | Structure of human GLUT3(N43T) bound to D-glucose.

a, b, Substrate specificity of GLUT3. The transport of D-[2-³H]glucose by recombinant GLUT3(N43T) protein was examined in proteoliposome-based counterflow assays in the presence of the indicated monosaccharides (**a**) and chemicals (**b**). All the experiments were independently repeated at least three times. Error bars represent s.d. L-DHA, L-dehydroascorbic acid; GlcN, glucosamine; GlcNAc, N-acetylglucosamine; G6P, glucose-6-phosphate. **c**, Overall structure of GLUT3 in the presence of D-glucose. GLUT3 is domain-coloured with pale green, pale cyan and pale yellow for the N-terminal, C-terminal and ICH domains, respectively. The 'omit' electron density of the bound ligand, shown as brown mesh, is contoured at 3σ. **d**, The structure of glucose-bound GLUT3 exhibits an outward-occluded conformation. The corresponding TMs in the four 3-TM repeats are coloured the same. The ligand is shown as black spheres. The surface electrostatic potential was calculated with PyMol⁴⁹. All structure figures were prepared with PyMol.

may be crucial for specific recognition by GLUT3. Supporting this notion, D-ribose, the C3 epimer of D-xylose, showed no inhibitory effect on the uptake of glucose (Fig. 1a and Extended Data Fig. 1a).

Among the chemicals tested, L-dehydroascorbic acid, which was suggested to be a substrate for GLUTs^{38,39}, effectively inhibited the counterflow of D-glucose by GLUT3. Addition of glucosamine, a reported substrate of GLUT2 (ref. 40), blocked up to 60% of the activity. In contrast, myo-inositol, N-acetylglucosamine or glucose-6-phosphate showed nearly no inhibition (Fig. 1b and Extended Data Fig. 1b). These results, consistent with characterizations performed in distinct systems⁴¹, verified the activity and selectivity of the glycosylation-site-eliminated recombinant GLUT3 protein.

Structure of GLUT3 in complex with D-glucose

The crystals of GLUT3 in the presence of D-glucose were obtained in the P2₁ space group using the lipidic cubic phase (LCP) method. The structure was solved with molecular replacement and refined to 1.5 Å

resolution (Fig. 1c, Extended Data Fig. 2a–d and Extended Data Table 1). The transmembrane region of GLUT3 exhibits a canonical major facilitator superfamily (MFS) fold with the 12 transmembrane segments (TMs) folded into the N-terminal and C-terminal domains, each comprising '3+3' inverted repeats⁴². TM7 and TM10 are discontinuous helices, hence named TM7a/7b and TM10a/10b, respectively (Fig. 1d). The N-terminal and C-terminal domains are connected by four helices (IC1–4), which together with the C-terminal helix IC5 constitute the intracellular helical (ICH) domain.

After assignment and refinement of the polypeptide chain of GLUT3 (residues 1–470), the omit electron density for a bound D-glucose was unambiguously recognized in the central cavity (Fig. 1c). The ligand is occluded from either side of the membrane in the outward-facing GLUT3. Therefore, the structure represents a substrate-bound and outward-occluded state (Fig. 1d). In addition to D-glucose, three monoolein molecules, which were the major components of LCP, were identified (Extended Data Fig. 2e, f).

Recognition of α- and β-anomers of D-glucose by GLUT3

Extensive inter-domain interactions constitute the extracellular and intracellular gates as well as the side walls that insulate the bound ligand from the surrounding environment (Extended Data Fig. 3 and Supplementary Discussion). The high resolution provides unprecedented clarity towards substrate recognition by GLUTs. It is of particular interest that the 2F_o – F_c electron density revealed the existence of both α- and β-D-glucose anomers. Despite the prevailing presence of β-D-glucose in aqueous solution, the α-anomer exhibits a dominant occupancy of approximately 69% in the refined structure (Fig. 2a, b).

D-glucose is bound asymmetrically within the central cavity (Fig. 2c, d and Extended Data Fig. 4a). Located halfway across the membrane height, it stands closer to the C-terminal domain, which provides the primary substrate accommodation site in the centre of its transport-path-facing surface. Six polar residues from the C-terminal domain, including Gln280 and Gln281 on TM7a, Asn286 on TM7b, Asn315 on TM8, Glu378 on TM10a, and Trp386 on TM10b, contribute eight hydrogen bonds to coordinate D-glucose. With respect to the N-terminal domain, D-glucose is closer to TM1 and TM5. Gln159 on TM5 is the only polar residue from the N-terminal domain that engages in hydrogen bonding to D-glucose (Fig. 2d–f).

The α- and β-D-glucose anomers are similarly coordinated by GLUT3 except for the variation at C1-OH (Fig. 2e, f). The α- and β-C1-OH are recognized by Trp386 and Gln280, respectively. The ring oxygen and C1-OH in both anomers are hydrogen bonded to the side group of Gln159. The C2, C3 and C6 hydroxyls are each coordinated by two hydrogen bonds, including those between C6-OH and Asn315 and Glu378, C2-OH and Gln280 and Trp386, and C3-OH and Gln281 and Asn286. The C4-OH forms a single hydrogen bond with Asn286 (Fig. 2e, f). In addition to the polar interactions, the carbon backbone of the sugar ring is surrounded by hydrophobic residues including Phe24 on TM1, Ile162 and Ile166 on TM5, Ile285 and Phe289 on TM7b, and Phe377 on TM10b (Fig. 2g). A monoolein molecule was found to participate in hydrogen bonding with D-glucose. Although the monoolein molecule should be physiologically irrelevant, similar indirect hydrogen bonds between glucose and polar residues may be mediated through water molecules under physiological conditions (Extended Data Fig. 4b).

Structures of GLUT3 bound to maltose

While working with crystallization of glucose-bound GLUT3, we attempted to capture the outward-open state of GLUTs by exploiting exofacial competitors^{43,44}. Details of this work can be found in the Supplementary Discussion and Extended Data Fig. 5. The structures of GLUT3 in complex with maltose were captured in two conformations, outward-open at 2.6 Å and outward-occluded at 2.4 Å resolutions, respectively (Fig. 3, Extended Data Fig. 6 and Extended Data

Figure 2 | Coordination of the α - and β -D-glucose anomers by GLUT3. **a, b**, Both α - and β -anomers of D-glucose are identified in the structure of GLUT3. The $2F_o - F_c$ electron density maps, shown as blue mesh, for α - and β -D-glucose are contoured at 2σ and 1σ , respectively. **c**, The bound glucose is predominantly coordinated by the C-terminal domain. A cut-open side view of the semi-transparent surface electrostatic potential is shown. **d**, The position of the bound glucose with respect to the C-terminal and N-terminal domains. The substrate-facing sides of the C-terminal and N-terminal domains are shown in surface and cartoon representations. Only the α -anomer is shown. **e, f**, Coordination of α - and β -D-glucose by GLUT3 through polar interactions. The α - (**e**) and β - (**f**) anomers are coloured black and silver, respectively. Hydrogen bonds are represented by red dashed lines. **g**, Coordination of glucose by GLUT3 through van der Waals interactions. The residues from the N-terminal and C-terminal domains are coloured green and cyan, respectively.

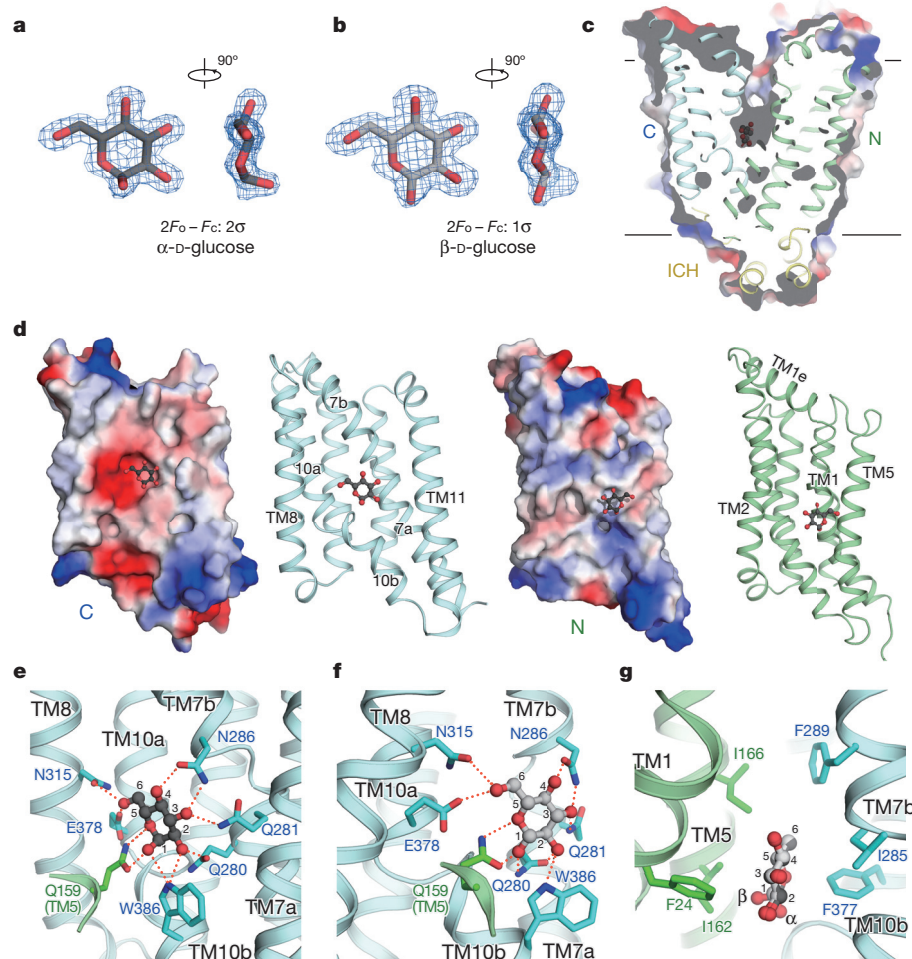


Table 2). Please refer to the Supplementary Discussion and Extended Data Fig. 7 for detailed analysis of maltose coordination in the two conformations. Briefly, the maltose-bound, outward-occluded structure is virtually indistinguishable from the glucose-bound GLUT3 with an overall root mean square deviation of 0.25 Å. The second glucose unit (Glc2) of maltose completely overlaps with D-glucose when the two structures are superimposed. The coordination of Glc2 is nearly identical to that of glucose except for Asn286, whose side group exhibits a distinct rotamer that recognizes the C3-OH of Glc1 (Extended Data Fig. 7a, b).

Comparison of the outward-open and -occluded states

Subtle relative rotation of the N-terminal and C-terminal domains towards the central cavity on the extracellular side is observed during the state transition from outward-open to outward-occluded. The N-terminal domains remain nearly rigid in the two structures, whereas prominent local rearrangements occur to TM7b in the C-terminal domain (Fig. 4a, b and Supplementary Video 1).

To achieve the conformational switch from outward-open to outward-occluded, TM7b, which consists of four helical turns, is partly unwound and bent in the middle such that the first two helical turns tilt inwards. Meanwhile these two helical turns undergo an axial rotation by approximately 60 degrees, leading to the relocation of the bulky residue Tyr290 and the substrate-coordinating residue Asn286 into the transport path (Fig. 4b). Notably, an invariant glycine (Gly284 in GLUT3) constitutes the kink preceding TM7b, which may provide the flexibility for the pronounced structural shift (Supplementary Fig. 1). Whereas the endofacial region of the TM domain remains nearly unchanged (Fig. 4a, right panel), the ICH domain exhibits minor intra-domain rearrangements with the inward

motion of helices IC1/2/3/5 by 1 to 2 Å during the outward-open to outward-occluded transition (Fig. 4c).

Alternating access

GLUT3 shares sequence identity of 66% and similarity of 80% with GLUT1 (Supplementary Fig. 1). Structural resolution of the two closely related GLUTs in three distinct conformations reveals the molecular basis for the alternating access cycle. The N-terminal domain remains almost rigid during a relative rotation of the two domains, while the discontinuous helices TM7/10 and their adjoining segments in the C-terminal domain undergo prominent local

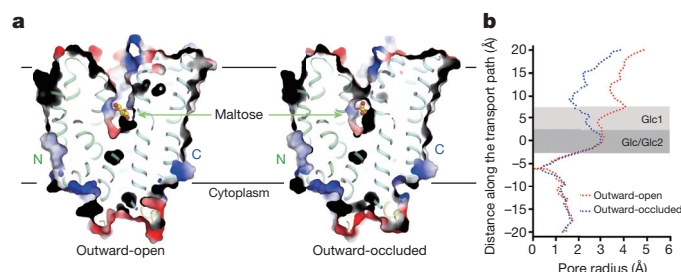
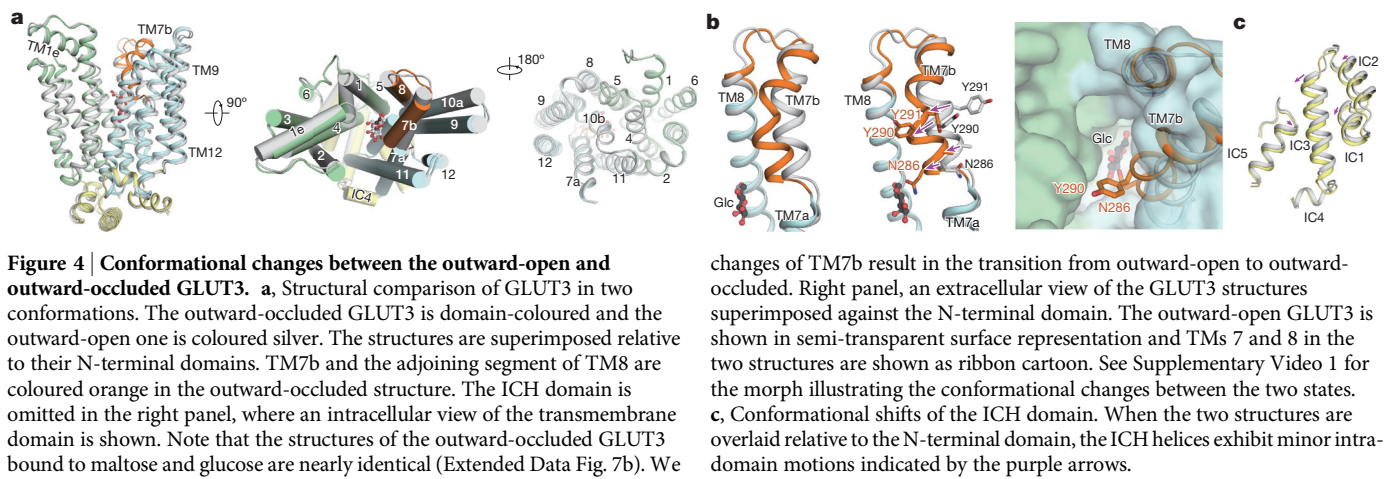


Figure 3 | Structures of maltose-bound GLUT3 in the outward-open and outward-occluded conformations. **a**, Structures of GLUT3 in complex with maltose in the outward-open and outward-occluded states. The cut-open views of the surface electrostatic potential are shown to compare the contours of the central cavities in the two structures. **b**, The van der Waals surface of the two GLUT3 structures was calculated with the program HOLE⁵⁰. The radii of the potential transport path are tabulated. The glucose and maltose binding zones are indicated by grey shades. Glc, glucose.



rearrangements. TM7b is gradually bent towards the central cavity, while TM10b swings away from the transport path during the outward to inward transition (Fig. 5a–c).

The rigidity of the N-terminal domain and the adaptability of the C-terminal domain may be determined by their distinct structural features (Extended Data Fig. 8). The internal core of the N-terminal domain is relatively hydrophilic. Seven water molecules were identified within the N-terminal domain of the GLUT3 structures. A continuous strip of hydrogen bonds connecting the water

changes of TM7b result in the transition from outward-open to outward-occluded. Right panel, an extracellular view of the GLUT3 structures superimposed against the N-terminal domain. The outward-open GLUT3 is shown in semi-transparent surface representation and TMs 7 and 8 in the two structures are shown as ribbon cartoon. See Supplementary Video 1 for the morph illustrating the conformational changes between the two states. **c**, Conformational shifts of the ICH domain. When the two structures are overlaid relative to the N-terminal domain, the ICH helices exhibit minor intra-domain motions indicated by the purple arrows.

molecules and the polar residues on TM1 and TM4 extends throughout the internal core of the N-terminal domain, likely supporting the rigidity of the N-terminal domain during conformational shift. In contrast, the C-terminal domain is highly hydrophobic. This ‘greasy’ interior may enable the structural adaptability of the C-terminal domain.

Consistent with the extensive polar interactions with the C-terminal domains, the bound D-glucose in the outward-facing GLUT3 and the glucoside of n-nonyl-β-D-glucopyranoside (β-NG)

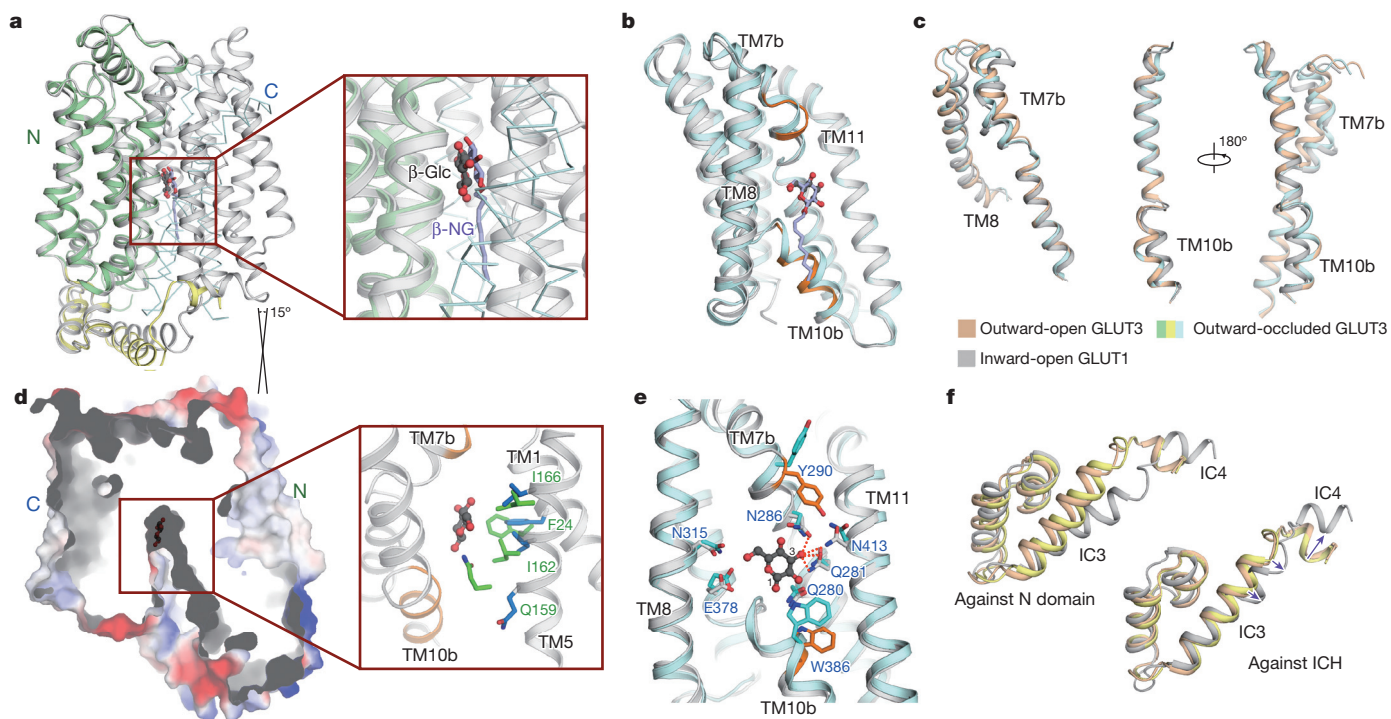


Figure 5 | Alternating access of the glucose-binding site in GLUTs.

a, Structural comparison of glucose-bound GLUT3 with the inward-open GLUT1. The two structures are superimposed relative to their N-terminal domains. GLUT3 is domain-coloured and GLUT1 (PDB accession code: 4PYP) is coloured silver. Inset, when the N-terminal domains of the two structures are superimposed, the bound ligands are not overlapped. For visual simplicity, the C-terminal domain of GLUT3 is shown as a thin ribbon. **b**, Intra-domain rearrangements of the C-terminal domains between the inward-open GLUT1 and outward-occluded GLUT3. The segments that exhibit local structural shifts are coloured orange in GLUT1. The structures of GLUT3 and GLUT1 are superimposed relative to their C-terminal domains in panels **b**–**e**. See Supplementary Video 2 for the morph that illustrates the conformational changes. **c**, TM7b and TM10b undergo prominent conformational changes in

the three indicated structures. **d**, Modelled glucose-binding in the inward-open GLUT1. The cut-open view of the semi-transparent surface electrostatic potential of GLUT1 is shown. Inset, the N-terminal domain residues may no longer engage in glucose-coordination in the inward-open conformation. The N-terminal domain residues that contribute to glucose binding in the outward-facing GLUT3 are shown as green sticks; the corresponding residues in the inward-open GLUT1 are coloured blue. All the residues are numbered as in GLUT3. **e**, Rearrangements of glucose-coordination between the outward-occluded GLUT3 and inward-open GLUT1. The C-domain ligand-binding residues that exhibit prominent changes between the two states are coloured orange in the inward-open conformation. **f**, Conformational changes of the ICH domain. The colour code of the three structures follows that in panel **c**.

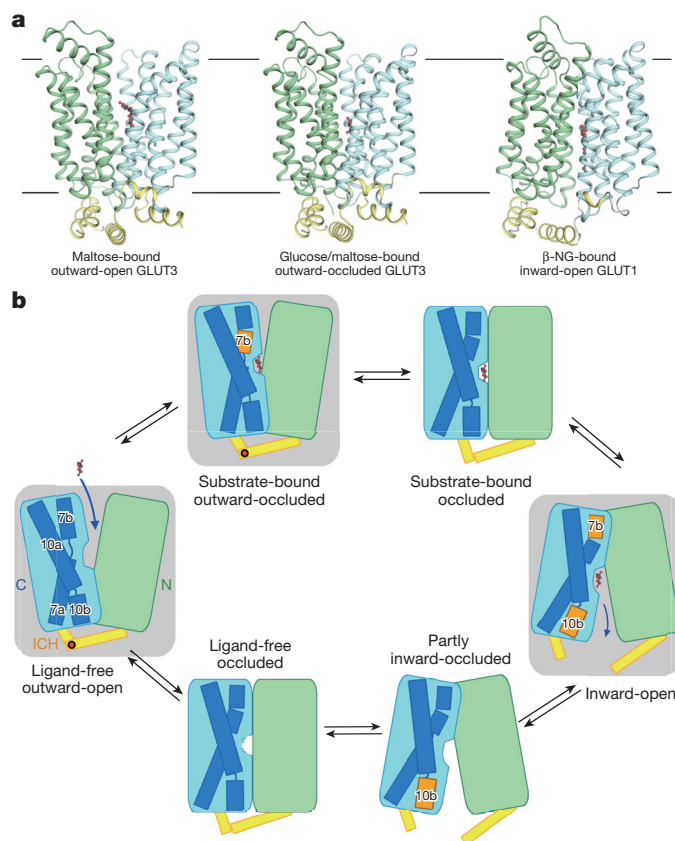


Figure 6 | Updated working model of GLUTs. **a**, An overview of the resolved structures of human GLUT1 and GLUT3. See Supplementary Video 3 for the morph that illustrates the conformational changes between the three structures. **b**, Schematic illustration of the alternating access cycle of GLUTs. The rearrangements of the substrate-binding site and the local structural shifts of TM7b and TM10b during the transport cycle are highlighted. The shaded states refer to the structures shown in panel **a**. The N-terminal, C-terminal and ICH domains are coloured green, cyan and yellow, respectively. The structural elements that undergo prominent local shifts during state transition are highlighted in orange.

in the inward-facing GLUT1 completely overlap when the C-terminal domains, but not N-terminal domains, are superimposed (Fig. 5a, b). We modelled β -D-glucose in the structure of inward-open GLUT1 at the position of the glucoside of β -NG (Fig. 5d). As all the concerned residues are invariant between GLUT1 and GLUT3, we used the residue numbering of GLUT3.

In the inward-open structure, the N-terminal domain residues no longer engage in ligand binding (Fig. 5d, inset). Within the C-terminal domain, local structural shifts result in rearrangement of the substrate-binding site (Fig. 5e). Trp386 loses contact with the ligand due to the outward swing of TM10b. Consequently, the ring oxygen and α -C1-OH may lose interactions with the protein in this conformation (Supplementary Video 2).

In addition to the local shift of TM7/TM10, the ICH domain also undergoes pronounced conformational changes. Because of the preferential association of helices IC1/2 with the N-terminal domain and IC4 with the C-terminal domain, the intervening helix IC3 may function like a 'door closer' that restrains the opening degree of the N-terminal and C-terminal domains on the intracellular side (Fig. 5f and Supplementary Video 3). Helix IC5 is missing in the structure of GLUT1, implying its intrinsic flexibility due to the potential loss of interactions with other ICH helices in the inward-open state³² (Extended Data Fig. 3c). The structural observation supports the notion that the ICH domain may serve as a latch to stabilize the outward-facing conformation of GLUTs^{32,33}.

Discussion

The structural analyses of GLUTs presented here and previously (Fig. 6a) provide the molecular basis to address the long-standing questions on substrate recognition and transport by GLUTs. The discrimination of α - and β -D-glucose by GLUTs and whether anomericization is required has been controversial for decades^{45–47}. The 1.5 Å structure of glucose-bound GLUT3 reveals that GLUTs can recognize both anomers, hence anomericization may not be required for D-glucose transport by GLUTs (Fig. 2e, f).

The structures also provide interpretation for the 'asymmetry' of ligand binding from the endo- and exofacial sides of GLUTs⁴⁸ (Figs 2, 5 and 6b). The C-terminal domain provides the primary yet partial substrate-binding site composed of polar residues including Gln280 and Gln281 on TM7b, Asn315 on TM8, and Glu378 on TM10a. Arrival of the substrate at the primary site from either the exo- or endofacial side may induce conformational changes including the local shifts of TM7b and TM10b along with the relative rotation of the N-terminal and C-terminal domains. Therefore, the substrate-binding site may undergo dynamic rearrangements during a transport cycle (Fig. 6b and Supplementary Video 2).

The advances in the structural elucidation of GLUT3 and GLUT1 provide the framework for understanding a wealth of experimental data of GLUTs over the last half century and provide the foundation for further kinetic and thermodynamic studies of this important class of transporters. The high-resolution structures offer important insight into substrate selectivity, which may illuminate the rational design and optimization of ligands targeting the Warburg effect through GLUTs.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 6 January; accepted 12 June 2015.

Published online 15 July 2015.

- Nelson, D. L. & Cox, M. M. *Lehninger Principles of Biochemistry* (W. H. Freeman, 2008).
- Hediger, M. A., Clemençon, B., Burrier, R. E. & Bruford, E. A. The ABCs of membrane transporters in health and disease (SLC series): introduction. *Mol. Aspects Med.* **34**, 95–107 (2013).
- Mueckler, M. & Thorens, B. The SLC2 (GLUT) family of membrane transporters. *Mol. Aspects Med.* **34**, 121–138 (2013).
- Manolescu, A. R., Witkowska, K., Kinnaird, A., Cessford, T. & Cheeseman, C. Facilitated hexose transporters: new perspectives on form and function. *Physiology* **22**, 234–240 (2007).
- LeFevre, P. G. Evidence of active transfer of certain non-electrolytes across the human red cell membrane. *J. Gen. Physiol.* **31**, 505–527 (1948).
- Kasahara, M. & Hinkle, P. C. Reconstitution and purification of the D-glucose transporter from human erythrocytes. *J. Biol. Chem.* **252**, 7384–7390 (1977).
- Mueckler, M. et al. Sequence and structure of a human glucose transporter. *Science* **229**, 941–945 (1985).
- Thorens, B. & Mueckler, M. Glucose transporters in the 21st century. *Am. J. Physiol. Endocrinol. Metab.* **298**, E141–E145 (2010).
- Dick, A. P., Harik, S. I., Klip, A. & Walker, D. M. Identification and characterization of the glucose transporter of the blood-brain barrier by cytochalasin B binding and immunological reactivity. *Proc. Natl Acad. Sci. USA* **81**, 7233–7237 (1984).
- Maher, F., Vannucci, S. J. & Simpson, I. A. Glucose transporter proteins in brain. *FASEB J.* **8**, 1003–1011 (1994).
- Thorens, B., Sarkar, H. K., Kaback, H. R. & Lodish, H. F. Cloning and functional expression in bacteria of a novel glucose transporter present in liver, intestine, kidney, and β -pancreatic islet cells. *Cell* **55**, 281–290 (1988).
- Fukumoto, H. et al. Sequence, tissue distribution, and chromosomal localization of mRNA encoding a human glucose transporter-like protein. *Proc. Natl Acad. Sci. USA* **85**, 5434–5438 (1988).
- Simpson, I. A. et al. The facilitative glucose transporter GLUT3: 20 years of distinction. *Am. J. Physiol. Endocrinol. Metab.* **295**, E242–E253 (2008).
- James, D. E., Brown, R., Navarro, J. & Pilch, P. F. Insulin-regulatable tissues express a unique insulin-sensitive glucose transport protein. *Nature* **333**, 183–185 (1988).
- Birnbaum, M. J. Identification of a novel gene encoding an insulin-responsive glucose transporter protein. *Cell* **57**, 305–315 (1989).
- Mueckler, M. Facilitative glucose transporters. *Eur. J. Biochem.* **219**, 713–725 (1994).
- Pascual, J. M. et al. GLUT1 deficiency and other glucose transporter diseases. *Eur. J. Endocrinol.* **150**, 627–633 (2004).
- Santer, R. et al. Mutations in GLUT2, the gene for the liver-type glucose transporter, in patients with Fanconi-Bickel syndrome. *Nature Genet.* **17**, 324–326 (1997).

19. Simpson, I. A., Chundu, K. R., Davies-Hill, T., Honer, W. G. & Davies, P. Decreased concentrations of GLUT1 and GLUT3 glucose transporters in the brains of patients with Alzheimer's disease. *Ann. Neurol.* **35**, 546–551 (1994).
20. Macheda, M. L., Rogers, S. & Best, J. D. Molecular and cellular regulation of glucose transporter (GLUT) proteins in cancer. *J. Cell. Physiol.* **202**, 654–662 (2005).
21. Amann, T. & Hellerbrand, C. GLUT1 as a therapeutic target in hepatocellular carcinoma. *Expert Opin. Ther. Targets* **13**, 1411–1427 (2009).
22. Shim, B. Y. *et al.* Glucose transporter 1 (GLUT1) of anaerobic glycolysis as predictive and prognostic values in neoadjuvant chemoradiotherapy and laparoscopic surgery for locally advanced rectal cancer. *Int. J. Colorectal Dis.* **28**, 375–383 (2013).
23. Ramani, P., Headford, A. & May, M. T. GLUT1 protein expression correlates with unfavourable histologic category and high risk in patients with neuroblastic tumours. *Virchows Arch.* **462**, 203–209 (2013).
24. Flavahan, W. A. *et al.* Brain tumor initiating cells adapt to restricted nutrition through preferential glucose uptake. *Nature Neurosci.* **16**, 1373–1382 (2013).
25. Airley, R. E. & Mobasher, A. Hypoxic regulation of glucose transport, anaerobic metabolism and angiogenesis in cancer: novel pathways and targets for anticancer therapeutics. *Chemotherapy* **53**, 233–256 (2007).
26. Kaira, K. *et al.* Biological significance of ¹⁸F-FDG uptake on PET in patients with non-small-cell lung cancer. *Lung Cancer* **83**, 197–204 (2013).
27. Gallamini, A., Zwarthoed, C. & Borra, A. Positron emission tomography (PET) in oncology. *Cancers* **6**, 1821–1889 (2014).
28. Calvaresi, E. C. & Hergenrother, P. J. Glucose conjugation for the specific targeting and treatment of cancer. *Chem. Sci.* **4**, 2319–2333 (2013).
29. Jardetzky, O. Simple allosteric model for membrane pumps. *Nature* **211**, 969–970 (1966).
30. Shi, Y. Common folds and transport mechanisms of secondary active transporters. *Annu. Rev. Biophys.* **42**, 51–72 (2013).
31. Smirnova, I., Kasho, V. & Kaback, H. R. Lactose permease and the alternating access mechanism. *Biochemistry* **50**, 9684–9693 (2011).
32. Deng, D. *et al.* Crystal structure of the human glucose transporter GLUT1. *Nature* **510**, 121–125 (2014).
33. Sun, L. *et al.* Crystal structure of a bacterial homologue of glucose transporters GLUT1–4. *Nature* **490**, 361–366 (2012).
34. Quistgaard, E. M., Low, C., Moberg, P., Tresaugues, L. & Nordlund, P. Structural basis for substrate transport in the GLUT-homology family of monosaccharide transporters. *Nature Struct. Mol. Biol.* **20**, 766–768 (2013).
35. Wisedchaisri, G., Park, M. S., Iadanza, M. G., Zheng, H. & Gonen, T. Proton-coupled sugar transport in the prototypical major facilitator superfamily protein XylE. *Nature Commun.* **5**, 4521 (2014).
36. Lacko, L., Wittke, B. & Geck, P. The temperature dependence of the exchange transport of glucose in human erythrocytes. *J. Cell. Physiol.* **82**, 213–218 (1973).
37. Chen, C. C. *et al.* Human erythrocyte glucose transporter: normal asymmetric orientation and function in liposomes. *Proc. Natl Acad. Sci. USA* **83**, 2652–2656 (1986).
38. Rumsey, S. C. *et al.* Glucose transporter isoforms GLUT1 and GLUT3 transport dehydroascorbic acid. *J. Biol. Chem.* **272**, 18982–18989 (1997).
39. Rumsey, S. C. *et al.* Dehydroascorbic acid transport by GLUT4 in *Xenopus* oocytes and isolated rat adipocytes. *J. Biol. Chem.* **275**, 28246–28253 (2000).
40. Uldry, M., Ibberson, M., Hosokawa, M. & Thorens, B. GLUT2 is a high affinity glucosamine transporter. *FEBS Lett.* **524**, 199–203 (2002).
41. Maher, F., Davies-Hill, T. M. & Simpson, I. A. Substrate specificity and kinetic parameters of GLUT3 in rat cerebellar granule neurons. *Biochem. J.* **315**, 827–831 (1996).
42. Yan, N. Structural advances for the major facilitator superfamily (MFS) transporters. *Trends Biochem. Sci.* **38**, 151–159 (2013).
43. Carruthers, A. & Helgerson, A. L. Inhibitions of sugar transport produced by ligands binding at opposite sides of the membrane. Evidence for simultaneous occupation of the carrier by maltose and cytochalasin B. *Biochemistry* **30**, 3907–3915 (1991).
44. Colville, C. A., Seatter, M. J., Jess, T. J., Gould, G. W. & Thomas, H. M. Kinetic analysis of the liver-type (GLUT2) and brain-type (GLUT3) glucose transporters in *Xenopus* oocytes: substrate specificities and effects of transport inhibitors. *Biochem. J.* **290**, 701–706 (1993).
45. Janoshazi, A. & Solomon, A. K. Initial steps of alpha- and beta-D-glucose binding to intact red cell membrane. *J. Membr. Biol.* **132**, 167–178 (1993).
46. Leitch, J. M. & Carruthers, A. α - and β -monosaccharide transport in human erythrocytes. *Am. J. Physiol. Cell Physiol.* **296**, C151–C161 (2009).
47. Carruthers, A. & Melchior, D. L. Transport of alpha- and beta-D-glucose by the intact human red cell. *Biochemistry* **24**, 4244–4250 (1985).
48. Barnett, J. E., Holman, G. D. & Munday, K. A. An explanation of the asymmetric binding of sugars to the human erythrocyte sugar-transport systems. *Biochem. J.* **135**, 539–541 (1973).
49. DeLano, W. L. The PyMOL molecular graphics system (Schrödinger, 2002).
50. Smart, O. S., Goodfellow, J. M. & Wallace, B. A. The pore dimensions of gramicidin A. *Biophys. J.* **65**, 2455–2460 (1993).

Supplementary Information is available in the online version of the paper.

Acknowledgements We thank the Tsinghua University Branch of China National Center for Protein Sciences (Beijing) for providing the facility support. The computation was completed on the "Explorer 100" cluster system of Tsinghua National Laboratory for Information Science and Technology. This work was supported by funds from the Ministry of Science and Technology of China (2015CB910101, 2011CB910501, 2014ZX09507003006) and National Natural Science Foundation of China (projects 31130002 and 31125009). The research of N. Y. was supported in part by an International Early Career Scientist grant from the Howard Hughes Medical Institute and an endowed professorship from Bayer Healthcare.

Author Contributions N.Y. conceived the project. D.D., P.S. and N.Y. designed all experiments. D.D., P.S., C.Y., X.J., L.X., W.R., K.H., M.Y. and S.F. performed the experiments. D.D., P.S., C.Y., M.K. and N.Y. analysed the data and contributed to manuscript preparation. N.Y. wrote the manuscript.

Author Information The X-ray crystallographic coordinates and structure factor files of the three human GLUT3 structures have been deposited in the Protein Data Bank (PDB) with the accession codes 4ZW9, 4ZWB, and 4ZWC. Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to N.Y. (nyan@tsinghua.edu.cn).

METHODS

Protein purification. The recombinant human glucose transporter GLUT1(N45T) was expressed and purified as described previously³². A similar protocol was applied for the expression and purification of human GLUT3 with minor modifications.

The codon-optimized cDNA of human GLUT3(N43T) was synthesized and subcloned into a modified PfstBac1 vector with an N-terminal 10 × His tag. Following the instructions for the Bac-to-Bac baculovirus expression system (Invitrogen), bacmids were generated in DH10Bac cells. The baculoviruses were generated and amplified in Sf-9 insect cells. For protein expression and purification, the Sf-9 cells were collected 72 h after viral infection, and disrupted using the dounce homogenizer for 80 cycles on ice. The membrane pellets were collected and homogenized in the buffer (25 mM MES pH 6.0 and 150 mM NaCl) containing protease inhibitors (aprotinin at 0.8 μM, pepstatin at 2 μM, and leupeptin at 5 μg ml⁻¹; Amresco), and then solubilized with 2% (w/v) *n*-dodecyl-β-D-maltoside (DDM, Anatrace) at 4 °C for 2 h. The insoluble fraction was precipitated by ultracentrifugation (150,000g) for 30 min at 4 °C. The supernatant was incubated with Ni-NTA resin (Qiagen) for an additional 30 min at 4 °C. The resin was then rinsed with the buffer containing 25 mM MES pH 6.0, 500 mM NaCl, 30 mM imidazole, and 0.06% (w/v) 6-cyclohexyl-1-hexyl-β-D-maltoside (CYMAL-6, Anatrace) three times. The protein was eluted with the washing buffer plus 300 mM imidazole and further concentrated to 10 mg ml⁻¹. For the transport assay, the protein was applied to Superdex 200 10/300 GL (SD200, GE Healthcare) pre-equilibrated with buffer containing 25 mM MES 6.0, 150 mM NaCl, and 0.06% (w/v) CYMAL-6. Peak fractions were collected for the transport assay. For crystallization, 50 mM D-glucose or maltose (Sigma) was added throughout the purification procedure. Concentrated protein was applied to HiTrap Desalting 5 ml (GE Healthcare) in the SD200 buffer plus 50 mM D-glucose or maltose.

Crystallization. GLUT3(N43T), which eluded our extensive crystallization trials using a conventional vapour-diffusion method, was crystallized with the lipid cubic phase (LCP) approach. To prepare the cubic phase for crystallization trials, the protein was concentrated to 30–40 mg ml⁻¹ before mixing with monoolein (Sigma) in 1:1.5 protein to lipid ratio (w/w) using a syringe lipid mixer⁵¹.

For crystallization of glucose-bound GLUT3, the 40 nl meso phase was mixed with 900 nl crystallization buffer for each condition on glass sandwich plates (Shanghai Fastal BioTech) using a robot arm Gryphon (ARI). Crystals appeared within one week with a typical size of 70 μm × 50 μm × 10 μm. These crystals diffracted X-rays to approximately 2.5 Å at SSRF beamline BL17U. Various strategies were exploited to optimize crystals. Finally, the mother liquor containing 28% (v/v) PEG400, 0.1 M HEPES pH 6.8, and 50 mM ammonium citrate gave rise to crystals with a size of approximately 140 μm × 100 μm × 20 μm.

For crystallization trials of maltose-bound GLUT3, 30–45 nl meso phase was overlaid with 800 nl of precipitant solution. Crystals appeared overnight under two similar conditions. The crystals corresponding to the outward-occluded conformation appeared in the crystallization conditions with 38–40% (v/v) PEG 400, 100 mM Mg(CHO₂)₂, 50 mM maltose, and 100 mM ADA pH 6.5. The crystals corresponding to the outward-open conformation appeared in 34% (v/v) PEG 400, 400 mM (NH₄)₂HPO₄, 50 mM maltose and 100 mM ADA pH 6.9. Crystals grew to a maximum size of about 20 × 20 × 5 μm³ at 20 °C within one week.

The crystals were collected using MicroMesh (M3-L18SP-50; MiTeGen) and immediately flash frozen in liquid nitrogen.

Data collection and processing. The data set of glucose-bound GLUT3 was collected at the microfocus beamline BL32XU at SPring-8, Japan. The LCP crystals of GLUT3 were fast-screened with the shutter free strategy developed by SPring-8. Employing the microfocus X-ray beam of size 1 μm × 12 μm and the 'helical collection strategy', we were able to collect a complete data set with a single elongated crystal. All diffraction data sets of maltose bound GLUT3 were also collected at the same beamline. There are generally 10–15 crystals in each MicroMesh. Diffraction data from 7 and 9 crystals in two different lattices were integrated and scaled using HKL2000, respectively⁵². Further processing was carried out with programs from the CCP4 suite⁵³. Data collection and structure refinement statistics are summarized in Extended Data Tables 1 and 2. The data sets were processed with the HKL2000 packages⁵².

Structure determination and refinement. The phase was solved by molecular replacement using PHASER⁵⁴ with GLUT1 (PDB code: 4PYP) as a searching model. The model was first modified by CHAINSAW⁵⁵, then the N-terminal domain and C-terminal domain were taken separately as input ensembles for PHASER. The model was further rebuilt in COOT⁵⁶ and refined with PHENIX⁵⁷. The sequence docking was further aided by sequence alignment with GLUT1. The percentages of occupancy for α- and β-D-glucose were determined by the PHENIX occupancy refinement strategy. The two anomers were considered

as a constrained occupancy group where the occupancies of atoms in α- and β-D-glucose are coupled. The calculated occupancies for the α- and β-anomers were 0.69 and 0.31, respectively. For Supplementary Video 1 the intermediate morphs were obtained with the multiple-chain morphing script^{58,59} for Crystallography & NMR system (CNS)^{60,61}. For Supplementary Video 2, a homology-based structural model of the inward-open GLUT3 was generated based on the structure of GLUT1 using the online SWISS-MODEL workspace^{62–64}.

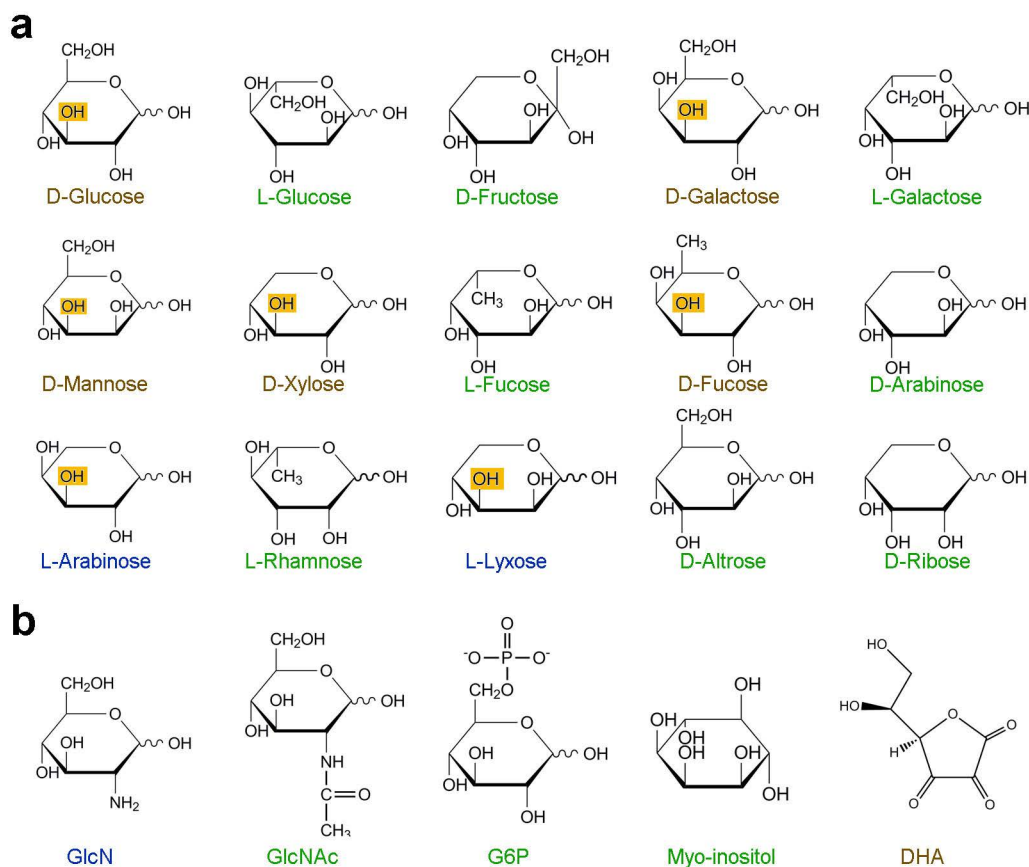
Preparation of liposomes and proteoliposomes. Liposomes were prepared as described previously³³. For the counterflow assay, proteoliposomes were prepared in solution containing KPM 6.5 buffer (50 mM potassium phosphate pH 6.5, 2 mM MgSO₄), 20 mg ml⁻¹ pre-extruded phospholipids (*E. coli* polar lipid extract; Avanti), and 50 mM D-glucose. After incubation with 1% *n*-octyl-β-D-glucopyranoside (β-OG; Anatrace) for 60 min at 4 °C, purified GLUT1(N45T) or GLUT3(N43T) (10 μg protein per mg lipid) was added and incubated for an additional 60 min at 4 °C. β-OG was removed by incubation with 240 mg ml⁻¹ Bio-Beads SM2 (Bio-Rad) overnight and an additional 2 h incubation with 120 mg ml⁻¹ Bio-Beads. Proteoliposomes were then frozen and thawed in liquid nitrogen for five cycles and extruded through membrane filter (PC Membranes 0.4 μm, Avanti). In order to remove the excessive glucose, proteoliposomes were collected by ultracentrifugation at 100,000g for 1 h and washed with ice-cold KPM 6.5 buffer. Finally, the proteoliposomes were resuspended in KPM 6.5 buffer to a final concentration of 100 mg ml⁻¹ (phospholipids).

Counterflow assay. For each assay, 2 μl of concentrated proteoliposomes (GLUT1 or GLUT3) prepared following the above protocol were added into 100 μl KPM 6.5 buffer plus 1 μCi D-[2-³H]glucose (specific radioactivity 21.5 Ci mmol⁻¹, PerkinElmer). The final concentration of the external D-[2-³H]glucose was 0.46 μM. The uptake of radiolabelled substrates was stopped at 30 s by rapidly filtering the solution through 0.22 μm filters (Millipore). The filter membranes were washed with 2 ml ice-cold KPM 6.5 buffer immediately, solubilized with 0.5 ml Optiphase HISAFE 3 (PerkinElmer) and used for liquid scintillation counting with MicroBeta JET (PerkinElmer). Liposomes without protein were tested as a negative control.

For competition assays, the indicated sugars and chemicals were added into the external KPM 6.5 at 50 mM and the transport reaction was stopped at 30 s. The reading of the competition assays was normalized against the one without competitor, which had the uptake of radiolabelled glucose set as 100%. All counter-flow assays were performed at 25 °C and repeated at least three times. Error bars represent s.d.

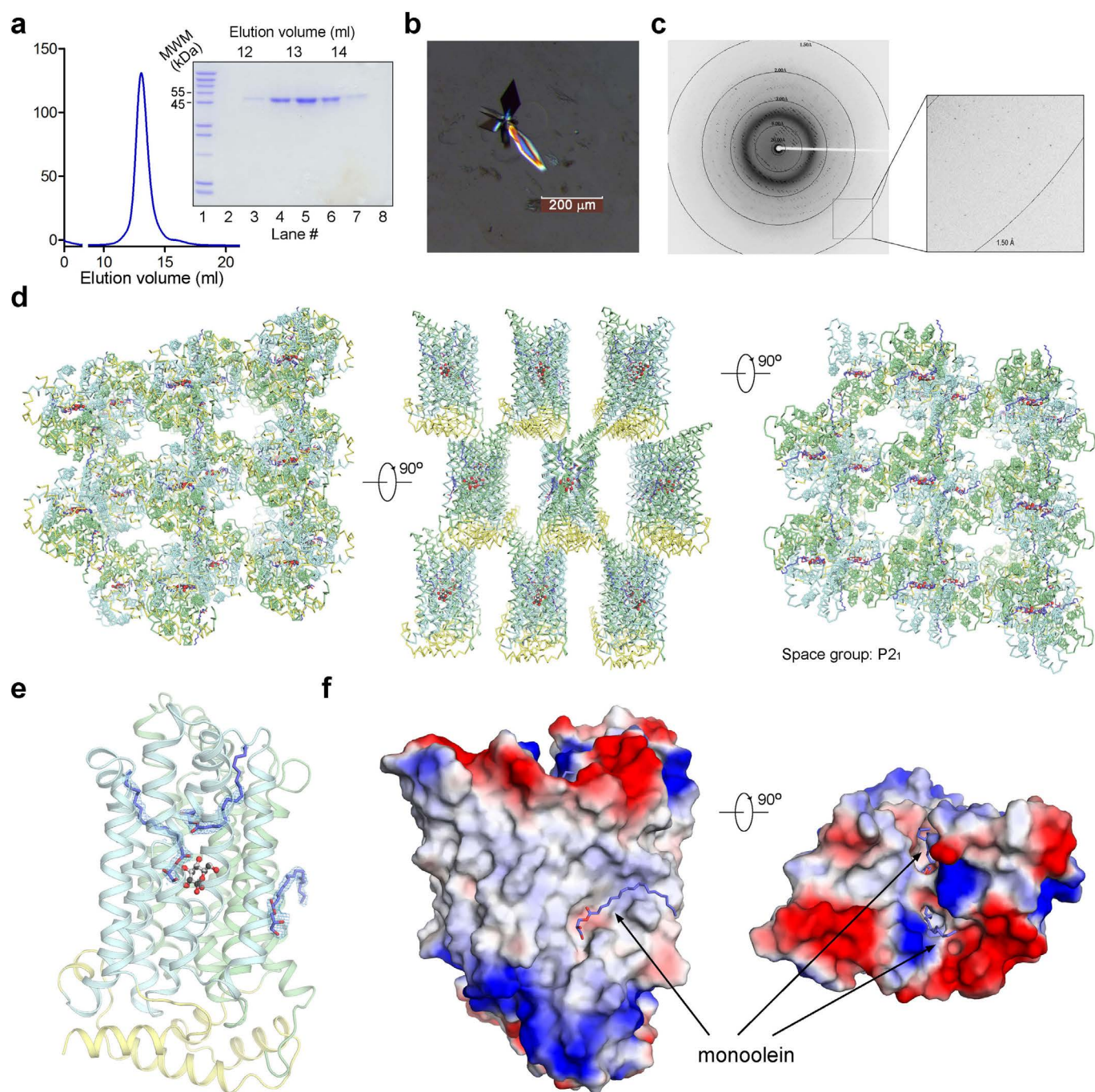
Data reporting. No statistical methods were used to predetermine sample size. The experiments were not randomized. The investigators were not blinded to allocation during experiments and outcome assessment.

- Caffrey, M. & Cherezov, V. Crystallizing membrane proteins using lipidic mesophases. *Nature Protocols* **4**, 706–731 (2009).
- Otwinski, Z. & Minor, W. Processing of X-ray diffraction data collected in oscillation mode. *Methods Enzymol.* **276**, 307–326 (1997).
- Collaborative Computational Project Number 4. The CCP4 suite: programs for protein crystallography. *Acta Crystallogr. D* **50**, 760–763 (1994).
- McCoy, A. J. *et al.* Phaser crystallographic software. *J. Appl. Crystallogr.* **40**, 658–674 (2007).
- Stein, N. CHAINSAW: a program for mutating pdb files used as templates in molecular replacement. *J. Appl. Crystallogr.* **41**, 641–643 (2008).
- Emsley, P. & Cowtan, K. Coot: model-building tools for molecular graphics. *Acta Crystallogr. D* **60**, 2126–2132 (2004).
- Adams, P. D. *et al.* PHENIX: building new software for automated crystallographic structure determination. *Acta Crystallogr. D Biol. Crystallogr.* **58**, 1948–1954 (2002).
- Echols, N., Milburn, D. & Gerstein, M. MolMovDB: analysis and visualization of conformational change and structural flexibility. *Nucleic Acids Res.* **31**, 478–482 (2003).
- Krebs, W. G. & Gerstein, M. The morph server: a standardized system for analyzing and visualizing macromolecular motions in a database framework. *Nucleic Acids Res.* **28**, 1665–1675 (2000).
- Brünger, A. T. *et al.* Crystallography & NMR system: a new software suite for macromolecular structure determination. *Acta Crystallogr. D* **54**, 905–921 (1998).
- Brünger, A. T. Version 1.2 of the Crystallography and NMR system. *Nature Protocols* **2**, 2728–2733 (2007).
- Arnold, K., Bordoli, L., Kopp, J. & Schwede, T. The SWISS-MODEL workspace: a web-based environment for protein structure homology modelling. *Bioinformatics* **22**, 195–201 (2006).
- Schwede, T., Kopp, J., Guex, N. & Peitsch, M. C. SWISS-MODEL: An automated protein homology-modeling server. *Nucleic Acids Res.* **31**, 3381–3385 (2003).
- Guex, N. & Peitsch, M. C. SWISS-MODEL and the Swiss-PdbViewer: an environment for comparative protein modeling. *Electrophoresis* **18**, 2714–2723 (1997).
- Thompson, J. D., Higgins, D. G. & Gibson, T. J. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* **22**, 4673–4680 (1994).



Extended Data Figure 1 | Chemical structures of the tested monosaccharides, glucose derivatives, and other chemicals in the competition assay. a, b, The monosaccharides (a) and chemicals (b) that exhibit potent inhibition to the

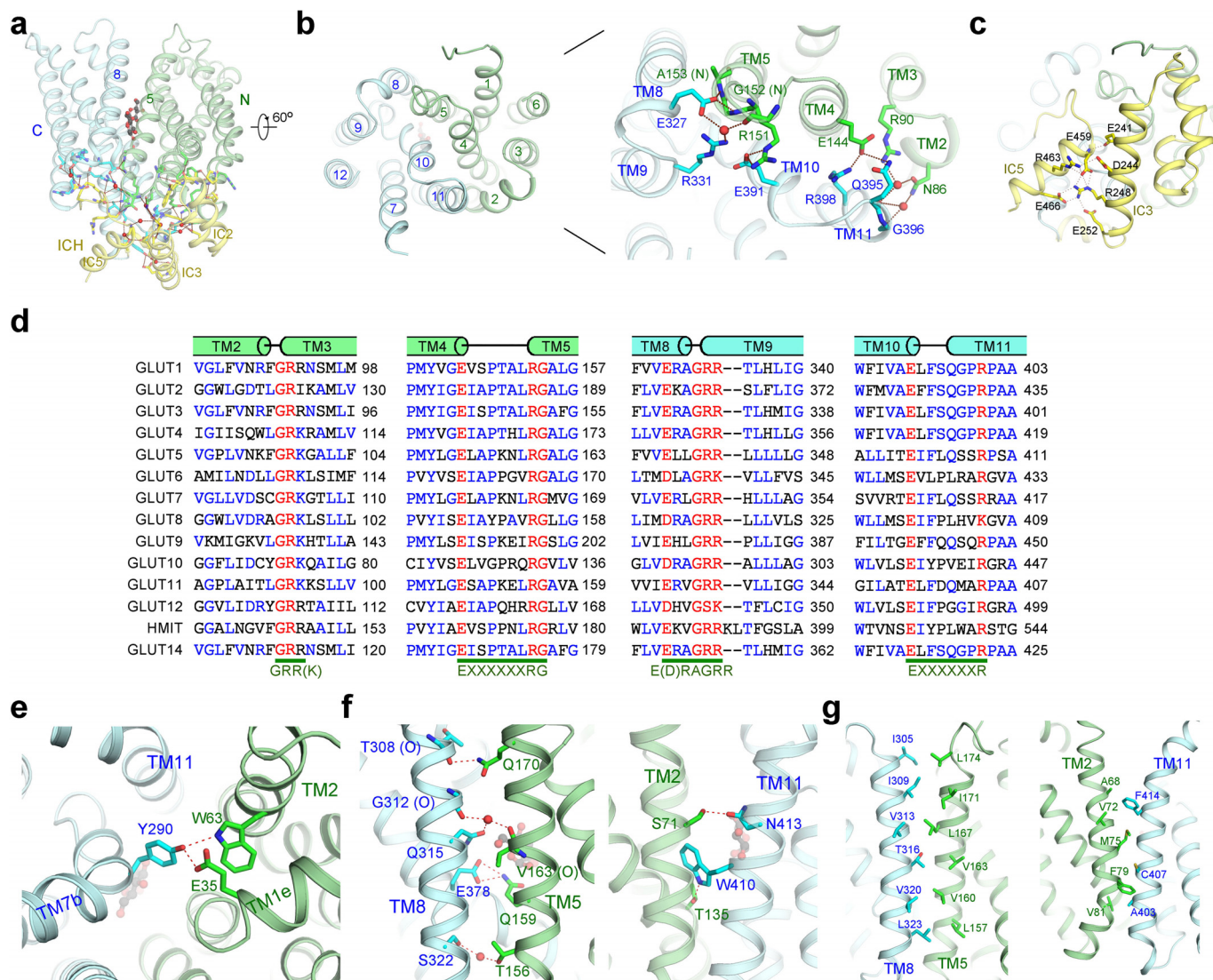
uptake of [2-³H]glucose in the proteoliposome-based counterflow assay are labelled brown. Those having moderate or no inhibition are labelled blue and green, respectively.



Extended Data Figure 2 | Protein purification, crystallization and structural determination of GLUT3(N43T) in the presence of D-glucose.

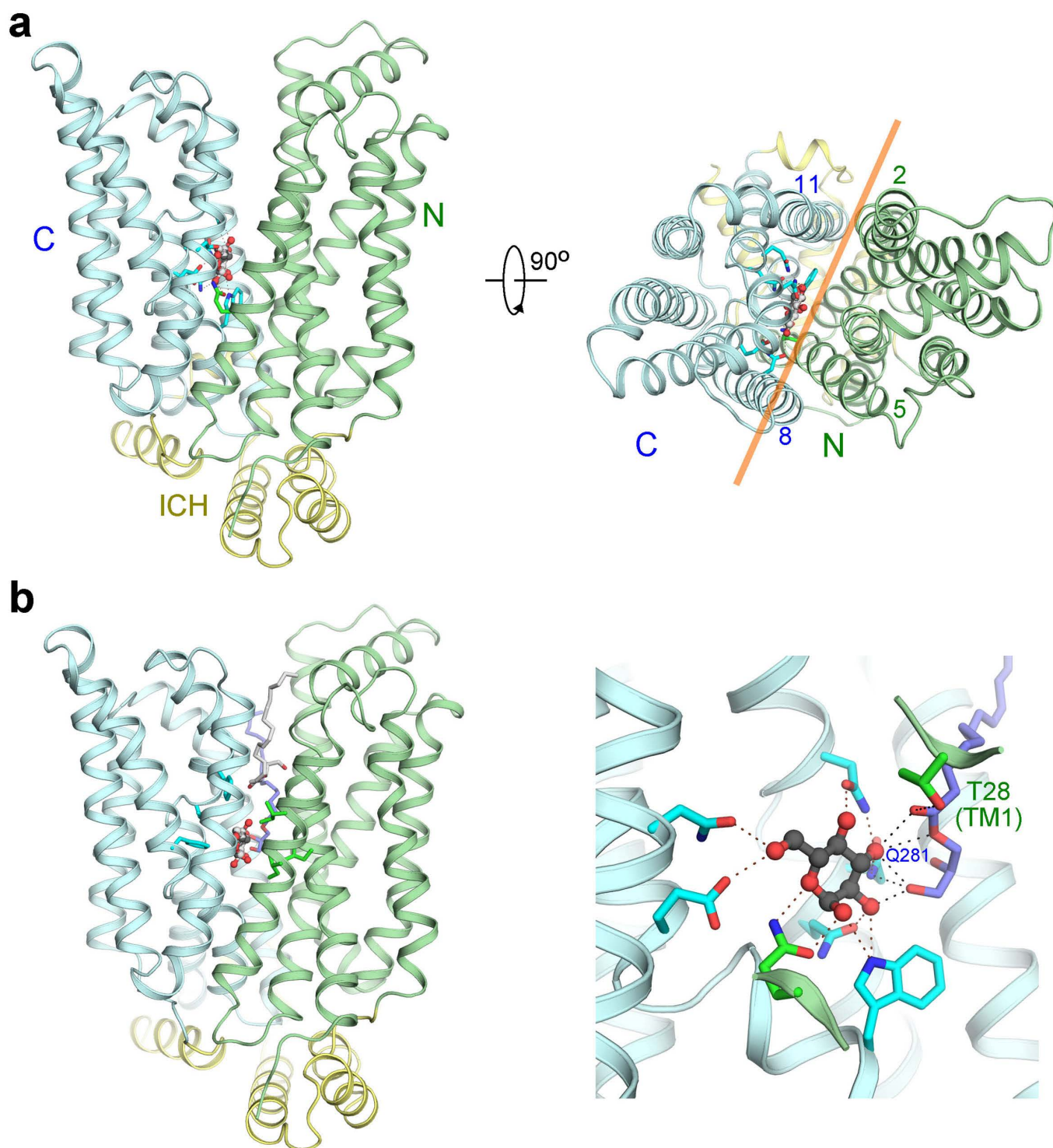
a, A representative chromatogram of the size-exclusion chromatography purification of GLUT3. The peak positions were applied for SDS-PAGE and followed by Coomassie blue staining. **b**, The crystal of GLUT3 used for X-ray diffraction which led to the final structural determination at 1.5 \AA resolution. **c**, A representative image of X-ray diffraction. The inset shows the resolution

limit beyond 1.5 \AA resolution. **d**, Crystal packing of GLUT3 in the space group of $P2_1$. Three perpendicular views are shown. Each GLUT3 is domain-coloured and shown as a ribbon. The bound glucose is shown as a black sphere and the bound lipid molecules are shown as blue sticks. **e**, Electron densities of three bound monoolein molecules. The $2F_o - F_c$ electron density maps for the bound lipid molecules are contoured at 1.0σ . **f**, Positions of the three bound monoolein molecules relative to the protein structure.



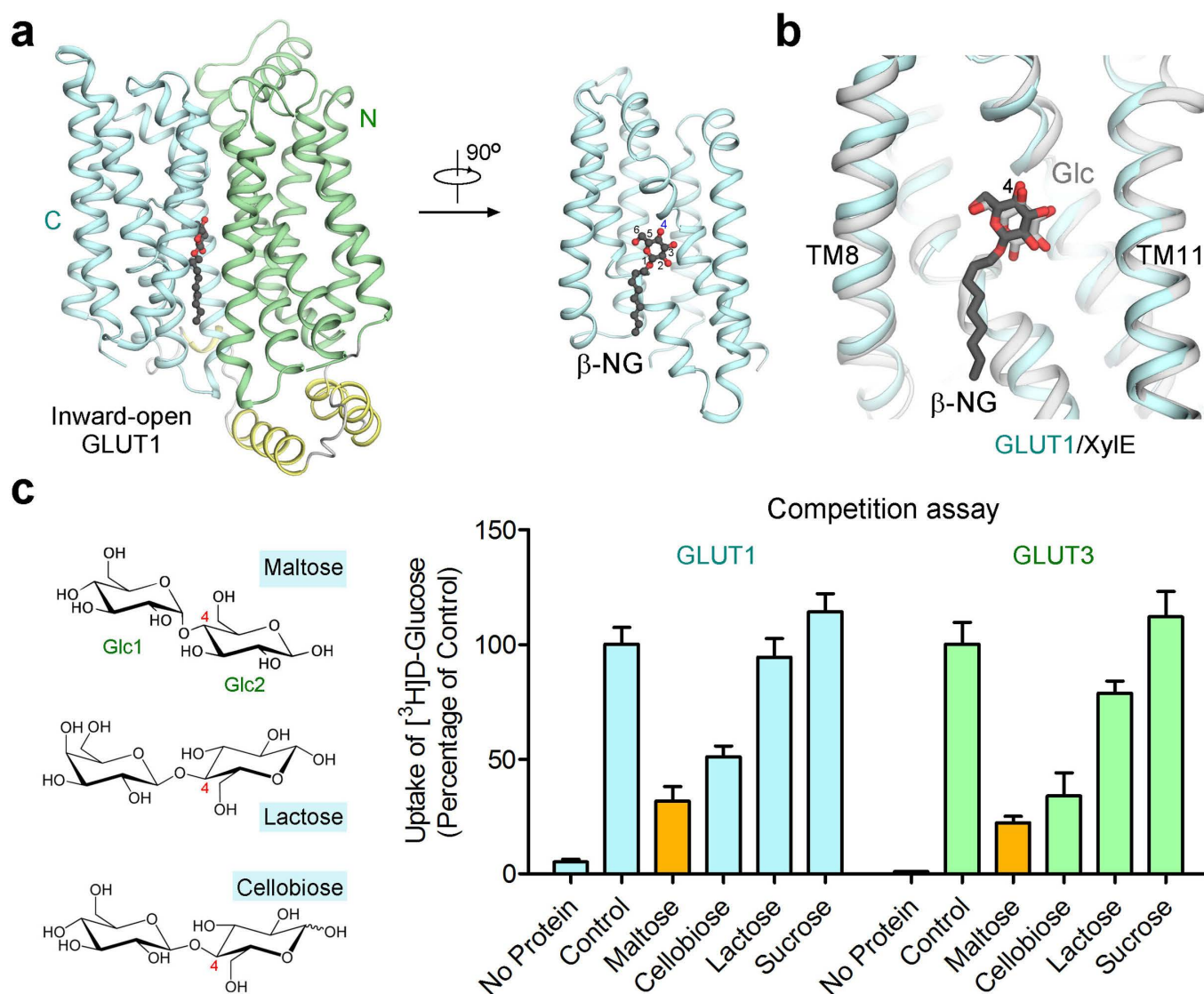
Extended Data Figure 3 | The intracellular and extracellular gates of GLUT3. **a**, The N-terminal, C-terminal and ICH domains interact with each other through an extensive network of direct and water-mediated hydrogen bonds on the intracellular side of the membrane. Water molecules are shown as red spheres. Hydrogen bonds are represented by red dashed lines. **b**, The polar interactions between the N-terminal and C-terminal domains on the intracellular side. An ICH domain-omitted intracellular view is shown. Note the pseudo two-fold symmetry of the overall structure and the interacting residues (inset on the right). These polar interactions partially constitute the intracellular gate of GLUT3 in the outward-facing conformation.

c, Intra-domain interaction of the ICH domain. The polar interactions between helices IC3 and IC5 are shown. **d**, The intracellular gate of GLUT3 involves the Sugar Porter family (SP)-signature motifs. Sequence alignment of the 14 human GLUTs were performed with ClustalW⁶⁵. The secondary structural elements and the SP motifs are shown on the top and bottom, respectively. The residues coloured red are the highly conserved residues that constitute the intracellular gate illustrated in panel **b**. **e**, The extracellular gate of GLUT3 in the outward-occluded state. An extracellular view is shown. **f**, **g**, The polar and hydrophobic residues mediating the lateral interactions of the N-terminal and C-terminal domains.



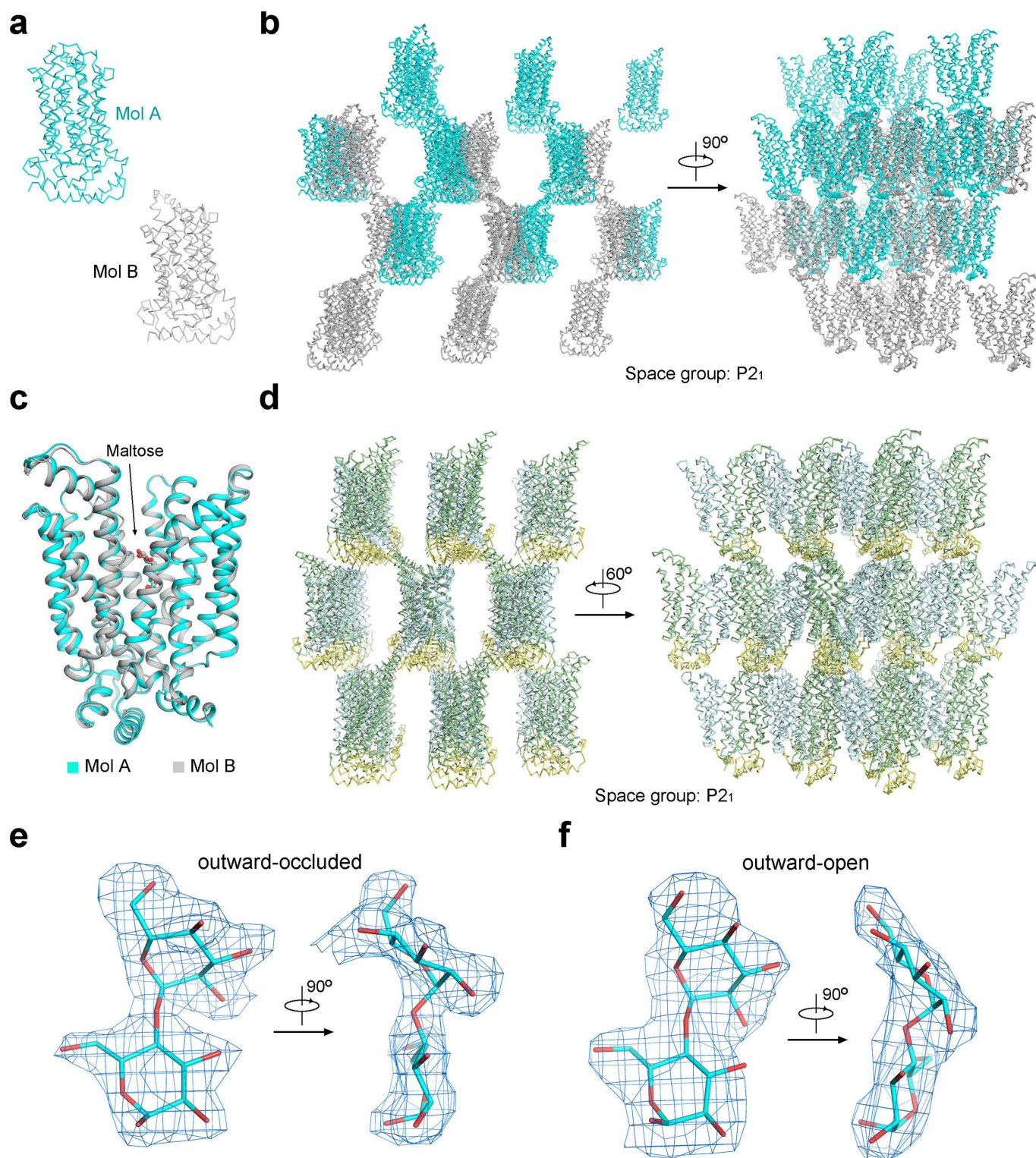
Extended Data Figure 4 | D-glucose coordination by GLUT3. **a**, The C-terminal domain provides the primary accommodation site for glucose in GLUT3. The α - and β -D-glucose anomers are coloured black and silver, respectively. The orange line in the right panel indicates the approximate interface between the N-terminal and C-terminal domains viewed from the extracellular side. Note that the ligand is located closer to the C-terminal

domain. **b**, One monoolein molecule contributes to substrate coordination. The two monoolein molecules bound in the cavity of GLUT3 are coloured silver and light purple. One monoolein molecule mediates indirect hydrogen bonds between C2 and C3 hydroxyl groups of the bound glucose with the side groups of Thr28 and Gln281 of GLUT3. The effect of monoolein on glucose binding and transport has not been characterized.



Extended Data Figure 5 | Attempts to obtain the structure of outward-facing GLUTs. **a**, Presence of the detergent molecule β -NG helped stabilize the inward-open conformation of GLUT1³². **b**, The glucopyranoside of β -NG and D-glucose are coordinated similarly by the inward-facing GLUT1 and the outward-facing XylE³³, respectively. In both structures, the C4-OH is positioned towards the extracellular side. **c**, Identification of potential ligands

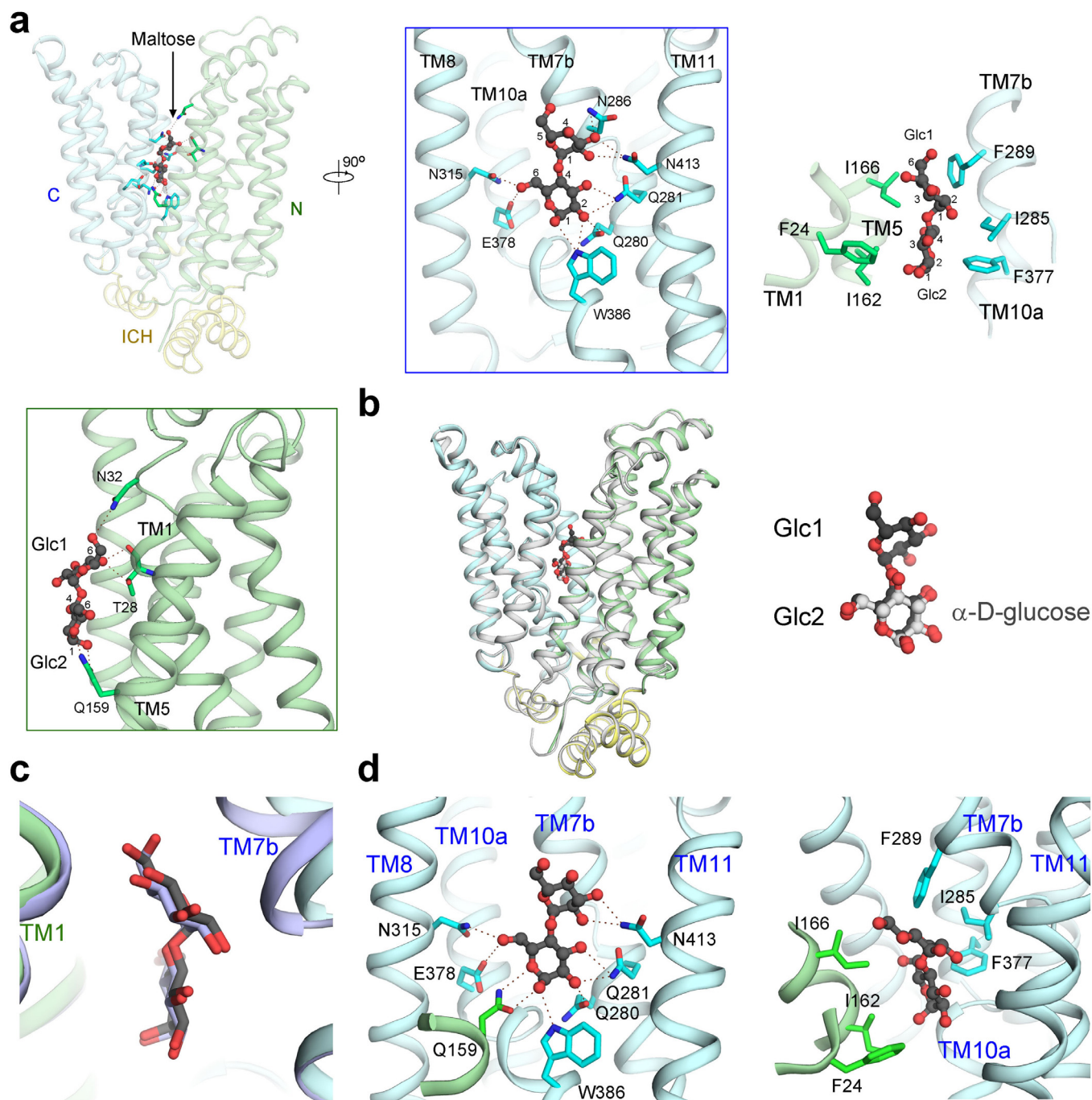
that may stabilize GLUTs in an outward-open conformation. The indicated disaccharides, where the C4-OH of D-glucose is condensed with another hexose, were tested for their abilities to inhibit glucose transport by GLUT1 or GLUT3 in the proteoliposome-based counterflow assay. Control refers to the conditions where no competitor was added.



Extended Data Figure 6 | Structure determination of GLUT3 in complex with maltose in the outward-open and outward-occluded conformations.

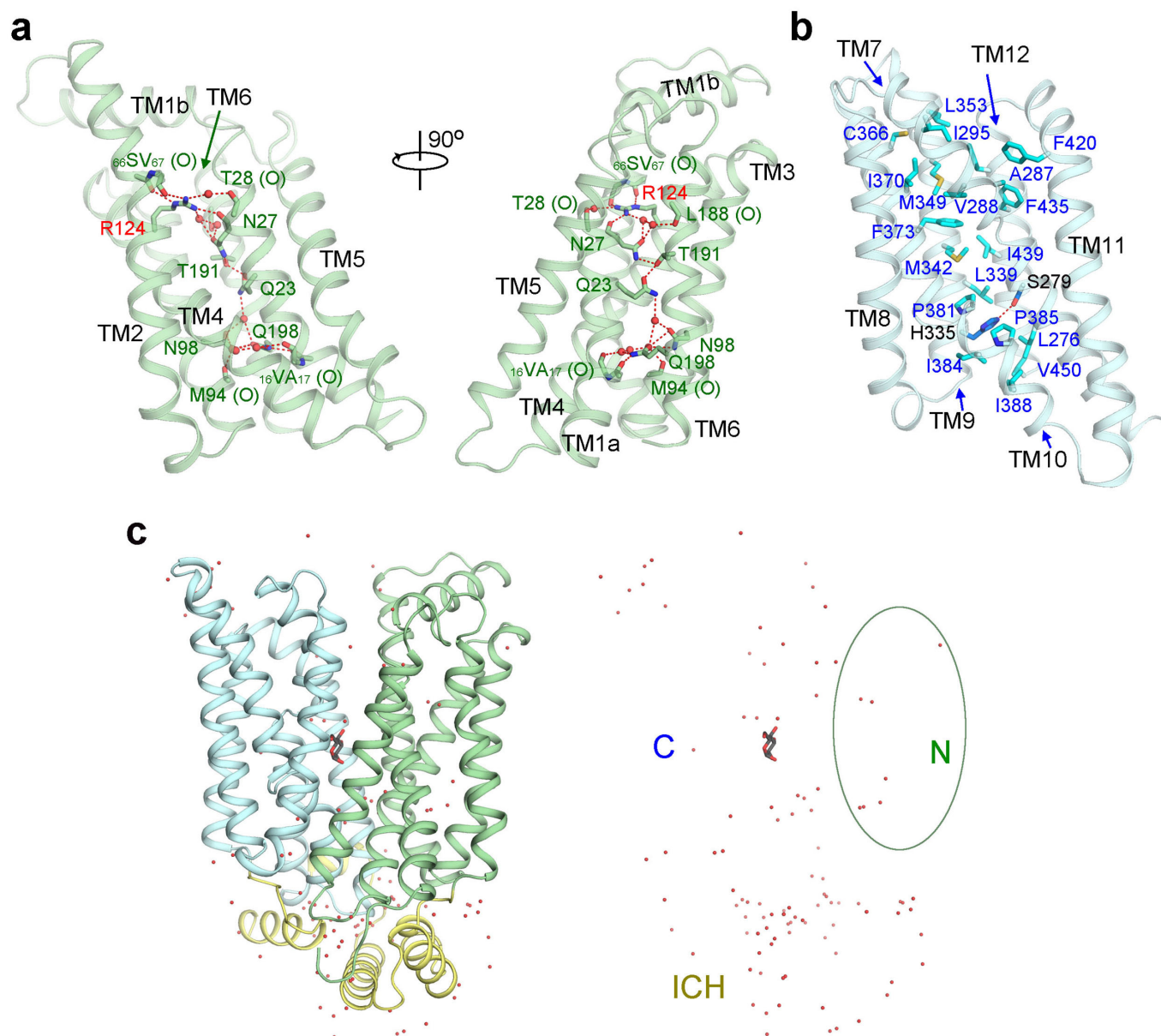
a, There are two GLUT3 molecules in each asymmetric unit of the outward-open structures. **b**, The crystal packing of the outward-open GLUT3. **c**, Structural superimposition of the two molecules in each asymmetric unit of the outward-open GLUT3. The two molecules exhibit nearly identical

conformations except for the extracellular loop regions. The focus was on molecule (Mol) A for structural analysis and comparison in the main text. **d**, The crystal packing of the outward-occluded GLUT3 bound to maltose. There is only one molecule in each asymmetric unit. **e**, **f**, The $2F_o - F_c$ electron density maps for the bound maltose in the outward-occluded (**e**) and outward-open (**f**) structures of GLUT3, both contoured at 1σ .



Extended Data Figure 7 | Maltose coordination in the outward-open and outward-occluded GLUT3. **a**, Coordination of maltose in the outward-occluded GLUT3 structure. Details of the polar interactions of maltose with residues from the N-terminal and C-terminal domains are shown in the insets on the bottom and right, respectively. Hydrogen bonds are represented by brown dashed lines. **b**, The structure of maltose-bound outward-occluded GLUT3 is nearly identical to that of the glucose-bound GLUT3. The second glucose unit (Glc2) of maltose completely overlaps with D-glucose, which is bound to the outward-occluded GLUT3. **c**, The bound maltose molecules are

positioned similarly in the outward-open and outward-occluded GLUT3. The two structures of GLUT3 are superimposed relative to the C-terminal domain. The outward-open GLUT3 is domain coloured and the outward-occluded structure is coloured pale purple. **d**, Maltose coordination in the outward-open GLUT3. Note that the coordination of Glc2 by the outward-open GLUT3 is identical to that by the outward-occluded GLUT3, while the polar residues in the N-terminal domain are not involved in the coordination of Glc1 in the outward-open structure.



Extended Data Figure 8 | The molecular basis underlying the rigidity and adaptability of the N-terminal and C-terminal domains, respectively.

a, Structural feature of the N-terminal domain of GLUT3. The interior of the N-terminal domain is held through a strip of hydrogen bonds, which may provide the molecular basis for the rigidity of the N-terminal domain during the alternating access cycle. **b**, The C-terminal domain of GLUT3 has a

hydrophobic core, which may allow the adaptability for the intra-domain shifts of the C-terminal domain during the transport cycle. **c**, Water distribution in the structure of glucose-bound GLUT3. The bound water molecules are shown as red spheres. Identical views are shown for the two panels except that the protein is omitted in the right panel to better illustrate the density of water molecules with respect to the N-terminal, C-terminal and ICH domains.

Extended Data Table 1 | Statistics of data collection and refinement of GLUT3 bound to D-glucose

Data	GLUT3/D-glucose
Integration Package	HKL2000
No. of crystals	1
Space Group	P2 ₁
Unit Cell (Å)	48.34, 118.13, 51.34
Unit Cell (°)	90, 102.67, 90
Wavelength (Å)	1.0000
Resolution (Å)	40~1.50 (1.55~1.50)
R _{merge} (%)	8.6 (57.8)
I/sigma	12.7 (3.6)
Completeness (%)	98.8 (91.8)
Number of measured reflections	444,895
Number of unique reflections	88,120
Redundancy	18.0 (2.0)
Wilson B factor (Å ²)	17.2
R _{work} / R _{free} (%)	17.53/19.42
No. atoms	
Protein	3,638
main chain	1,897
side chain	1,741
Substrate	12
Water	101
Others	75
Average B value (Å ²)	
Protein	25.08
main chain	22.93
side chain	27.43
Substrate	18.91
Water	40.50
Others	48.40
R.m.s. deviations	
Bonds (Å)	0.008
Angle (°)	1.095
Ramachandran plot statistics (%)	
Most favorable	97.0
Additionally allowed	3.0
Generously allowed	0.0
Disallowed	0.0

Extended Data Table 2 | Statistics of data collection and refinement of GLUT3 bound to maltose.

Data	Outward-occluded GLUT3/maltose	Outward-open GLUT3/maltose
Integration Package	HKL2000	HKL2000
No. of crystals	7	9
Space Group	P2 ₁	P2 ₁
Unit Cell (Å)	48.46, 119.49, 53.91	78.11, 121.88, 96.08
Unit Cell (°)	90, 103.75, 90	90, 108.14, 90
Wavelength (Å)	1.0000	1.0000
Resolution (Å)	40~2.40 (2.49~2.40)	40~2.60 (2.69~2.60)
R _{merge} (%)	13.0 (35.1)	14.3 (51.8)
I/sigma	12.7 (3.6)	10.5 (2.6)
Completeness (%)	94.6 (91.4)	95.0 (97.3)
Number of measured reflections	72,217	175,711
Number of unique reflections	22,422	49,602
Redundancy	3.2 (2.8)	3.5 (3.6)
Wilson B factor (Å ²)	33.3	84.6
R _{work} / R _{free} (%)	18.17 / 22.50	21.54 / 25.08
No. atoms		
Protein	3,593	7,186
main chain	1,872	3,744
side chain	1,721	3,442
Ligand	23	46
Water	36	52
Others	-	125
Average B value (Å ²)		
Protein	37.29	50.00
main chain	35.63	49.32
side chain	39.08	50.73
Ligand	38.24	59.00
Water	33.65	42.16
Others	-	60.8
R.m.s. deviations		
Bonds (Å)	0.009	0.010
Angle (°)	1.116	1.260
Ramachandran plot statistics (%)		
Most favourable	96.8	96.4
Additionally allowed	3.0	3.6
Generously allowed	0.2	0.0
Disallowed	0.0	0.0

Structure and mechanism of the mammalian fructose transporter GLUT5

Norimichi Nomura^{1,2,3*}, Grégory Verdon^{4,5,6*}, Hae Joo Kang^{4,5,6*}, Tatsuro Shimamura^{1,2,3}, Yayoi Nomura^{1,2,3}, Yo Sonoda⁴, Saba Abdul Hussien⁷, Aziz Abdul Qureshi⁷, Mathieu Coincon⁷, Yumi Sato^{1,3}, Hitomi Abe¹, Yoshiko Nakada-Nakura^{1,3}, Tomoya Hino^{1,2}, Takatoshi Arakawa^{1,2}, Osamu Kusano-Arai⁸, Hiroko Iwanari⁸, Takeshi Murata^{1,2,3,9}, Takuya Kobayashi^{1,2,3}, Takao Hamakubo⁸, Michihiro Kasahara¹⁰, So Iwata^{1,2,3,4,5,6,9} & David Drew^{4,7}

The altered activity of the fructose transporter GLUT5, an isoform of the facilitated-diffusion glucose transporter family, has been linked to disorders such as type 2 diabetes and obesity. GLUT5 is also overexpressed in certain tumour cells, and inhibitors are potential drugs for these conditions. Here we describe the crystal structures of GLUT5 from *Rattus norvegicus* and *Bos taurus* in open outward- and open inward-facing conformations, respectively. GLUT5 has a major facilitator superfamily fold like other homologous monosaccharide transporters. On the basis of a comparison of the inward-facing structures of GLUT5 and human GLUT1, a ubiquitous glucose transporter, we show that a single point mutation is enough to switch the substrate-binding preference of GLUT5 from fructose to glucose. A comparison of the substrate-free structures of GLUT5 with occluded substrate-bound structures of *Escherichia coli* Xyle suggests that, in addition to global rocker-switch-like re-orientation of the bundles, local asymmetric rearrangements of carboxy-terminal transmembrane bundle helices TM7 and TM10 underlie a ‘gated-pore’ transport mechanism in such monosaccharide transporters.

GLUT transporters belong to the solute carrier 2 family (*SLC2*) and, so far, 14 different isoforms (GLUT1–GLUT14) have been identified^{1,2}. Except for GLUT13, GLUT transporters are uniporters, which facilitate the diffusion of monosaccharides like glucose and fructose across the cell membrane in a concentration-dependent manner^{1,2}. Each GLUT transporter shows a distinct pattern of tissue distribution, gene regulation, substrate preference and kinetic properties^{1,2}. For example, GLUT1 is distributed in a wide range of tissues, including the blood–brain barrier, and is essential for glucose transport into the brain³, whereas GLUT4 is mostly localized to skeletal muscles and adipose tissue, and is the major insulin-stimulated glucose transporter^{1,4}. GLUT5 is the only member specific to fructose^{5,6}, and together with GLUT2, which transports fructose in addition to glucose, they make up the major fructose transporters in the body^{7,8}. GLUT5 is primarily expressed in the small intestine, but lower levels are also expressed in brain, adipose tissue, kidney, testes and skeletal muscle^{5,9,10}. GLUT activity is also associated with various diseases^{1,11}. For instance, increased GLUT5 expression has been linked to several metabolic disorders^{12,13} and several types of cancers such as breast cancer¹⁴, because of the higher energy demand of cancer cells stimulating sugar uptake, the so-called Warburg effect¹⁵. GLUT transporters belong to the larger major facilitator superfamily (MFS), the members of which have a fold consisting of two symmetrical six transmembrane helix (TM) bundles^{16,17}. Within the MFS they belong to a subfamily of sugar porter transporters^{18,19}, whose members are found in all domains of life and are important targets for industrial and biomedical applications²⁰. Recently, an open inward-facing

structure of human GLUT1 was reported with a bound sugar from a detergent head-group in the substrate-binding site and compared to previous structures of the related *Escherichia coli* D-xylose:H⁺ symporter Xyle in the outward- and inward-occluded conformations, suggesting a ‘rocker-switch’-type transport mechanism^{21–23}. However, as little is known about the molecular basis of substrate binding and release in GLUT transporters, their alternating-access mechanism is yet to be fully understood.

Open outward and inward GLUT5 structures

Rat and bovine GLUT5 (rGLUT5 and bGLUT5, respectively) that share ~81% sequence identity to human GLUT5 were selected and optimized for structural studies using fluorescence-based screening methods (Methods). rGLUT5 was crystallized in complex with an Fv antibody fragment (rGLUT5–Fv; see Methods). The rGLUT5–Fv and bGLUT5 structures were solved by molecular replacement (MR) and refined against data extending up to 3.3 Å and 3.2/4.0 Å (anisotropic data), respectively (Extended Data Tables 1 and 2, Extended Data Fig. 1 and Methods). The GLUT5 structure shows the typical MFS fold, plus five additional helices on the intracellular side, one at the C terminus (ICH5) and the other four, ICH1–ICH4, located between the amino- and C-terminal TM bundles (Fig. 1). bGLUT5 crystallized in an open inward-facing conformation (Fig. 1), and although human GLUT1 (hGLUT1) and bGLUT5 share only 43% sequence identity, their inward-facing structures superimpose well, with a root mean squared deviation (r.m.s.d.) of 1.12 Å for 364 pairs of Cα atoms (Methods and Extended Data Fig. 2a). The rGLUT5–Fv structure

¹Department of Cell Biology, Graduate School of Medicine, Kyoto University, Konoe-cho, Yoshida, Sakyo-ku, Kyoto 606-8501, Japan. ²Japan Science and Technology Agency, ERATO, Iwata Human Receptor Crystallography Project, Yoshida-Konoe-cho, Sakyo-ku, Kyoto 606-8501, Japan. ³Japan Science and Technology Agency, Research Acceleration Program, Membrane Protein Crystallography Project, Yoshida-Konoe-cho, Sakyo-ku, Kyoto 606-8501, Japan. ⁴Division of Molecular Biosciences, Imperial College London, London SW7 2AZ, UK. ⁵Membrane Protein Laboratory, Diamond Light Source, Harwell Science and Innovation Campus, Didcot, Chilton, Oxfordshire OX11 0DE, UK. ⁶Research Complex at Harwell, Rutherford Appleton Laboratory, Harwell, Oxford, Didcot, Oxfordshire OX11 0FA, UK. ⁷Centre for Biomembrane Research, Department of Biochemistry and Biophysics, Stockholm University, SE-106 91 Stockholm, Sweden. ⁸Department of Quantitative Biology and Medicine, Research Center for Advanced Science and Technology, University of Tokyo, 4-6-1 Komaba, Meguro-ku, Tokyo 153-8904, Japan. ⁹Systems and Structural Biology Center, RIKEN, 1-7-22 Suehiro-cho, Tsurumi-ku, Yokohama 230-0045, Japan. ¹⁰Laboratory of Biophysics, School of Medicine, Teikyo University, Hachioji, Tokyo 192-0395, Japan.

*These authors contributed equally to this work.

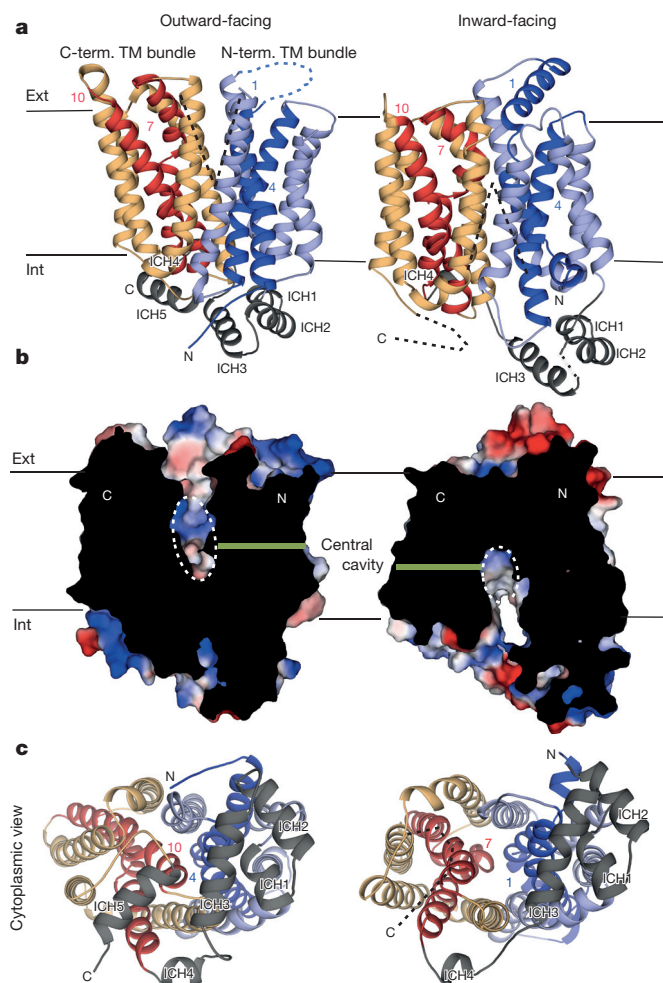


Figure 1 | Structures of rat GLUT5 in the open outward-facing conformation and bovine GLUT5 in the open inward-facing conformation. **a**, Ribbon representation of open outward-facing rat GLUT5 (left) and open inward-facing bovine GLUT5 (right) structures, viewed in the plane of the membrane. TM1 and TM4, and TM2, TM3, TM5 and TM6 in the N-terminal TM bundle are coloured in blue and light-blue, respectively. TM7 and TM10, and TM8, TM9, TM11 and TM12 in the C-terminal TM bundle are coloured in red and yellow–brown, respectively. The intracellular domain helices ICH1 to ICH5 are shown in grey. **b**, Slab through the surface electrostatic potential of the open outward-facing (left) and open inward-facing (right) GLUT5 structures, as viewed within the plane of membrane, which highlight the accessibility of the sugar to the central cavity (shown as a dotted ellipse). **c**, Ribbon diagrams of GLUT5 viewed from the cytoplasm in the open outward-facing (left) and inward-facing (right) conformations.

shows an open outward-facing conformation, which is a state that has not been observed previously in any of the related sugar porter structures^{22–25} (Fig. 1). The open outward-facing conformation is possibly stabilized by the Fv fragment, which binds to the ICHs (Extended Data Fig. 3).

Central fructose-binding site of GLUT5

The GLUT5 substrate-binding site is closely related to those of hGLUT1 and *E. coli* XylE (refs 21, 22) (Fig. 2a and Extended Data Fig. 2b). Many of the residues lining the central cavity are conserved between GLUT5 and hGLUT1, and include Ile169, Ile173, Gln166, Gln287, Gln288, Asn324 and Trp419 (Fig. 2a and Extended Data Fig. 4). In GLUT5, Trp419 is the only tryptophan positioned in the substrate-binding site (Fig. 2a and Extended Data Fig. 5a), and it is essential for transport²⁶. Consistent with rGLUT5 transport activity (Extended Data Fig. 6a), strong quenching of tryptophan fluorescence

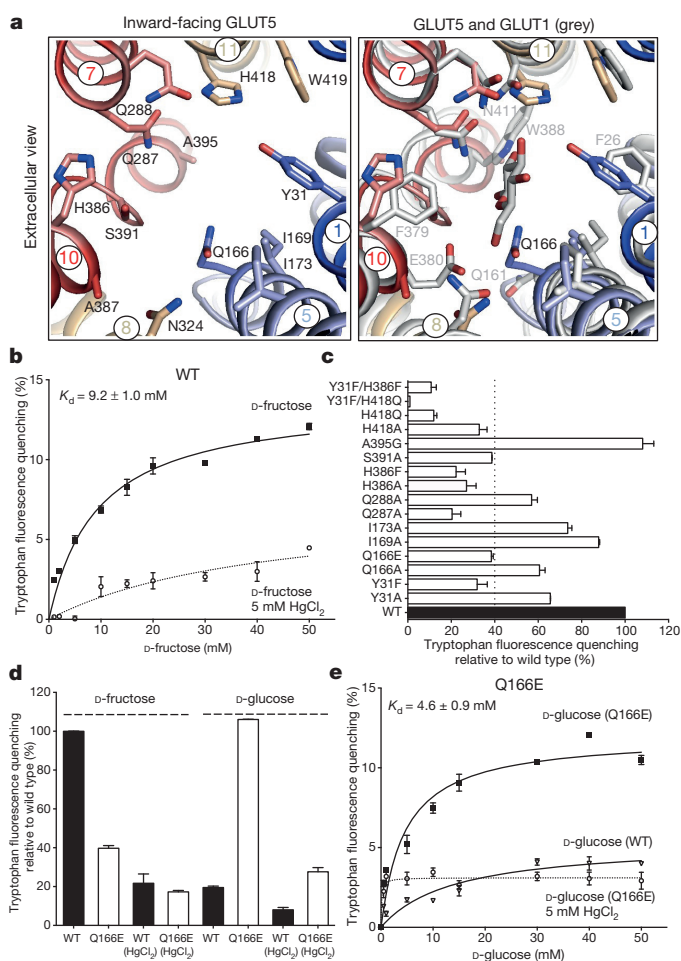


Figure 2 | The fructose-binding site of GLUT5. **a**, The substrate-binding site in the inward-facing bGLUT5 structure (left panel; coloured as in Fig. 1) is very similar to the inward-facing hGLUT1 structure (right panel; light grey). To facilitate comparison to rGLUT5, bGLUT5 residues are labelled with rGLUT5 numbering. For hGLUT1, only Q161 and all other residues that are different in bGLUT5 are labelled. The D-glucopyranoside moiety of bound *n*-nonyl-β-D-glucopyranoside in hGLUT1 is shown in stick representation. **b**, D-fructose binding to GLUT5 as measured by tryptophan (Trp) fluorescence quenching (excitation 295 nm; emission 338 nm) after addition of increasing concentrations of D-fructose to wild-type (WT; black squares) and to wild-type protein that had been previously incubated with the GLUT inhibitor $HgCl_2$ (open circles). **c**, Trp fluorescence quenching for purified substrate-binding site mutants after addition of 40 mM D-fructose (open bars) relative to wild type (filled bar). **d**, Trp fluorescence quenching after addition of either 40 mM D-fructose or D-glucose to purified wild type (filled bars) or Q166E (open bars); pre-incubation with the inhibitor $HgCl_2$ is indicated. **e**, Trp fluorescence quenching after addition of increasing concentrations of D-glucose to either purified Q166E (filled squares), Q166E previously incubated with $HgCl_2$ (open circles) or wild type (open triangles). In all experiments errors bars indicate s.e.m.; $n = 3$.

could be observed with the addition of D-fructose, but not with the addition of L-fructose or known GLUT1 substrates like D-glucose, D-galactose or D-mannose (Extended Data Fig. 5). Using this assay, the affinity of rGLUT5 for D-fructose was measured to have a dissociation constant (K_d) around 6–9 mM (Fig. 2b and Extended Data Fig. 6b), which is similar to that reported for human GLUT5 (refs 5, 8). Throughout this study tryptophan fluorescence quenching of rGLUT5 was henceforth used to assess substrate binding.

Gln166 is the only conserved residue in the substrate-binding site that is positioned differently between the hGLUT1 and GLUT5 structures (Fig. 2a). In hGLUT1 the equivalent glutamine (Gln161) points away, rather than towards, the central cavity as it has probably

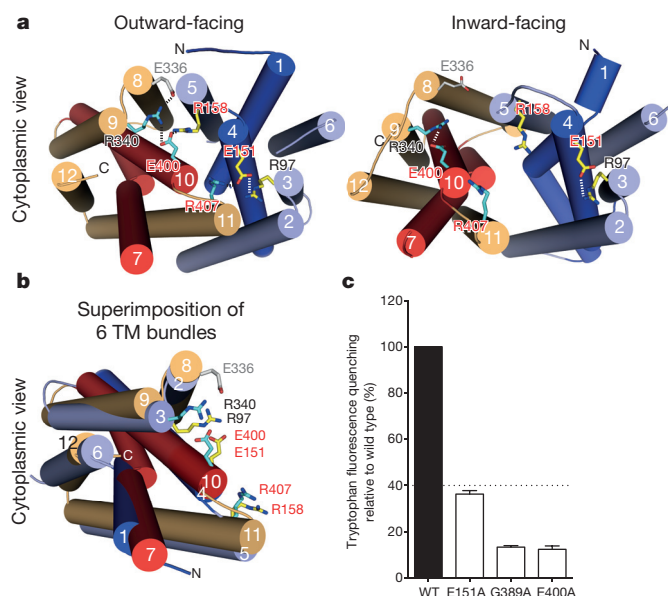


Figure 3 | Inter-TM salt bridges form between bundle cytoplasmic ends in the outward-facing conformation. **a**, Cartoon representation of GLUT5 as viewed from the cytoplasm in the outward-facing (left) and inward-facing (right) conformations. ICHs are not shown for clarity. TMs are coloured as in Fig. 1a, and residues forming salt bridges are shown as sticks. To facilitate comparison to rGLUT5, bGLUT5 residues are labelled with rGLUT5 numbering. Breakage of conserved salt bridges as seen here for GLUT5 has also been predicted for Xyle²³. **b**, Superimposition of the N- and C-terminal 6-TM bundles. Strictly conserved and pseudo-symmetry-related charged residues forming the salt bridges are labelled and are shown as sticks. **c**, D-fructose binding as measured by Trp fluorescence quenching after incubation with 40 mM D-fructose for purified wild-type GLUT5 (filled bar) and single alanine mutations of key acidic residues E400 and E151 that form inter-bundle salt bridges, and G389, which is located in the hinge point of TM10 critical for TM10 conformational change (open bars). Trp fluorescence quenching by D-fructose for these mutants is displayed as a percentage of wild-type binding. In all experiments errors bars indicate s.e.m.; $n = 3$.

re-orientated to accommodate the acyl chain of the bound detergent in the sugar-binding site (Fig. 2a). In Xyle the equivalent glutamine (Gln168) is orientated similarly as in GLUT5 (Extended Data Fig. 2b). Other residues that line the binding site, but are not conserved between GLUT5 and hGLUT1, are Tyr31 (Phe26 in hGLUT1), His386 (Phe379), Ala395 (Trp388), His418 (Asn411) and Ser391 (Gly384) (Fig. 2a and Extended Data Fig. 4). Single alanine mutants of each of these residues show strongly reduced D-fructose binding (Fig. 2c). The alanine mutants that show the weakest D-fructose binding are Tyr31 (TM1), Gln287 (TM7), His386 (TM10) and His418 (TM11). Except for Tyr31, these residues belong to the C-terminal TM bundle, indicating an asymmetrical binding mode of the sugar between the N- and C-terminal TM bundles in GLUT5, as seen in the sugar-bound hGLUT1 and Xyle structures^{21,22}. The substrate-binding cavity is deeper in GLUT5, because the tryptophan residue found at the bottom of the cavity in hGLUT1 (Trp388) and Xyle (Trp392) is replaced with alanine (Ala395) in GLUT5 (Fig. 2a and Extended Data Fig. 6c). Trp388 in hGLUT1 is critical for inhibition by cytochalasin B²⁶, whereas GLUT5 is insensitive to this inhibitor⁵.

GLUT7 is the closest isoform to GLUT5 and transports both D-fructose and D-glucose²⁷. Among the substrate-binding-site residues of GLUT5 almost all are identical to those in GLUT7 (Extended Data Fig. 4). The most notable difference is that Gln166 in GLUT5 is replaced in GLUT7 with glutamate (Extended Data Fig. 4). In rGLUT5, the substitution of Gln166 (TM5) with glutamate clearly weakens D-fructose binding (Fig. 2c, d). Furthermore, the Q166E mutant now shows robust binding to D-glucose with an apparent affinity (K_d) of ~4 mM (Fig. 2d, e). The introduced carboxylate is

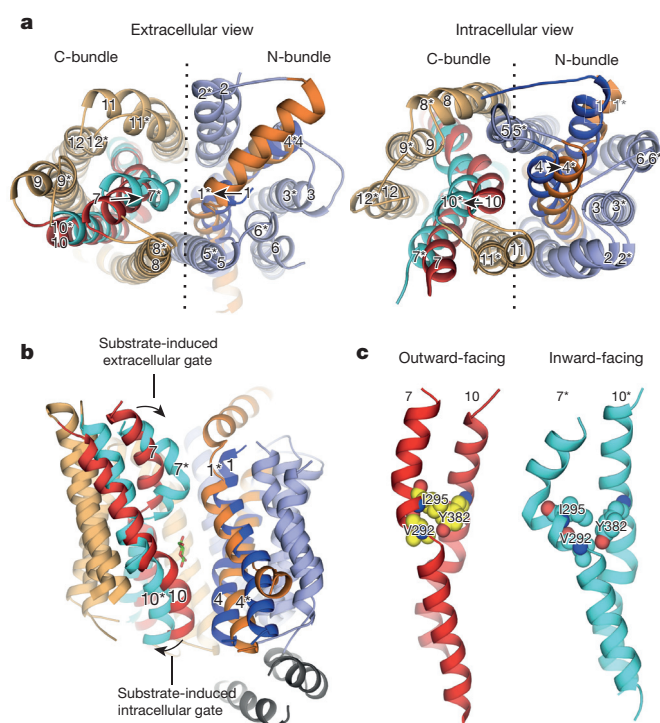


Figure 4 | Substrate-induced gates are predominantly formed by TM7 and TM10 in the C-terminal bundle. **a**, Superposition of GLUT5 open outward- and inward-facing (*) structures, as viewed from the extracellular (left) and intracellular (right) side of the membrane. TMs are coloured as in Fig. 1a, except inward-facing helices TM1*, TM4* shown in orange and TM7*, TM10* shown in cyan. ICHs have been removed for clarity. **b**, Superimposition of the GLUT5 open outward- and inward-facing structures as viewed in the plane of the membrane. For clarity, TM5, TM5*, TM8 and TM8* are not shown. Cavity-closing contacts are mostly formed by TM1* and TM7* on the extracellular side in the inward-facing conformation and by TM4 and TM10 on the intracellular side in the outward-facing conformation. These TMs are the first TMs in each of the four 3-TM repeats of the MFS fold^{16,37}. D-xylose, as it is in the occluded-outward-facing Xyle structure (PDB 4GBY), is shown in stick representation. With the inward movement of TM7, conserved tyrosine residues are likely to occlude the substrate from exiting, as seen for the equivalently located tyrosine residues in the substrate-occluded Xyle structure²² and as supported by D-fructose binding data (Extended Data Fig. 6d). The opening movement of TM10 to enable cytosolic substrate release has been described previously for Xyle^{23,24} and other unrelated MFS transporters^{38,39}. **c**, Interactions between hydrophobic residues between TM7 and TM10 in the outward-facing conformation (left) are lost in the inward-facing conformation (right). To facilitate comparison to rGLUT5, bGLUT5 residues are labelled with rGLUT5 numbering.

likely to have a similar role to Glu380 (TM10) in hGLUT1 (Fig. 2a), which is essential for D-glucose transport²⁸.

Salt bridges stabilize the outward conformation

A common feature observed in recent structures of MFS transporters are inter-TM bundle salt bridges that break and form near the central cavity during transport^{16,29,30}. In H^+ -coupled transporters, it is thought that the protonation state of these salt bridges is coupled to substrate binding and further to structural transitions^{16,29,30}. In GLUT5, no salt bridges are observed near the central cavity in either conformation. Instead, inter-TM bundle salt bridges are observed only in the outward-facing conformation and far from the central cavity, linking the cytoplasmic ends of TM3, TM4 and TM5 in the N-terminal TM bundle to those of TM9, TM10 and TM11 in the C-terminal TM bundle (Fig. 3a). Specifically, Glu151 (TM4) forms a salt bridge to both Arg97 (TM3) and Arg407 (TM11), and similarly Glu400 (TM10) forms a salt bridge to both Arg158 (TM5) and Arg340

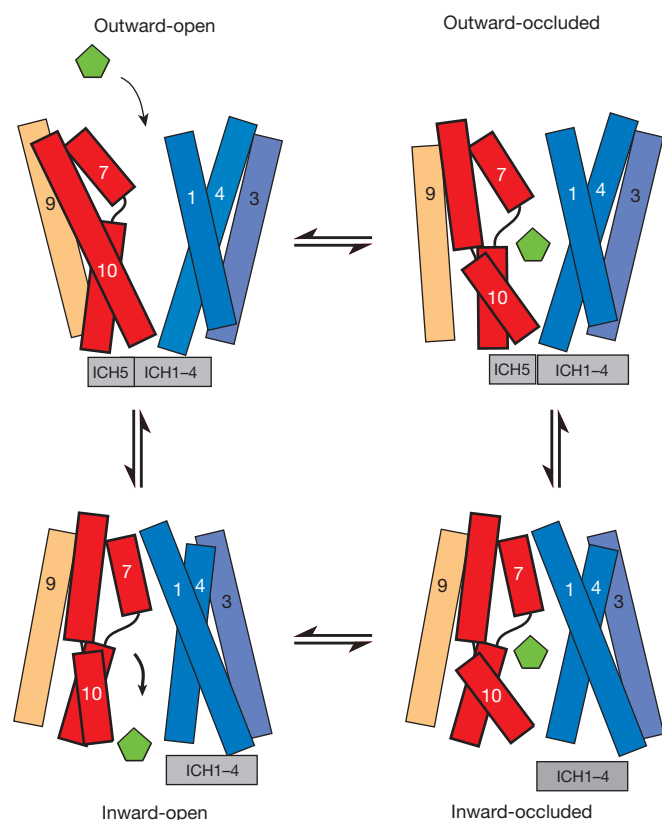


Figure 5 | Alternating-access transport mechanism in GLUT5. Schematic representation of the rocker-switch-type movement of the N- and C-terminal TM bundles and of the local, gating conformational changes of TM7 and TM10 supporting a gated-pore-type transport mechanism in GLUT5.

(TM9). Owing to their strict conservation these residues make up the well-described sugar porter signatures^{18,19} (Extended Data Fig. 4), and are related by a pseudo-two-fold symmetry axis running through the centre of the transporter and perpendicular to the membrane plane (Figs 1a and 3b). In the inward-facing conformation the inter-TM bundle salt-bridge pairs Glu400–Arg158 and Glu151–Arg407 are separated by some ~ 7 Å and 13 Å, respectively (Fig. 3a). In the inward-facing conformation, no inter-TM bundle salt bridges are formed on the extracellular side, indicating that in the absence of substrate GLUT proteins may favour the outward-facing conformation. Consistently, the substitution of Glu400, Arg407 and Glu336 equivalent residues in GLUT4 to neutral amino acids arrests the transporter in an inward-facing conformation³¹. Moreover, in the outward-facing conformation Glu336 (TM8) in the C-terminal TM bundle forms a salt bridge to Arg340, which is connected to the inter-TM bundle Glu400–Arg158 salt bridge (Fig. 3a). The equivalent glutamate to Glu336 in hGLUT1 (Glu329) was substituted to glutamine to stabilize the inward-facing hGLUT1 structure²¹. The salt bridges formed between inter-TM bundles seem to be coupled also to substrate binding in GLUT5, because substituting Glu151 and Glu400 with alanine shows significantly reduced D-fructose binding (Fig. 3c).

The ICH domain is below this cytoplasmic salt-bridge network, and positioned similarly, with respect to the N- and C-terminal TM bundles, in both outward-open GLUT5 and outward-occluded Xyle structures²² (Extended Data Fig. 7a). ICH1, ICH2 and ICH3 are linked together by salt bridges, as previously shown for Xyle and GLUT1 (refs 21, 22). In GLUT5, ICH5 lacks charged residues that could interact with ICH1–ICH3 (Extended Data Fig. 7b), as it does in the outward-occluded Xyle structure²². Rather, the N-terminal helices ICH1–ICH3 interact with the C-terminal bundle via a salt bridge formed between Glu252 in ICH3 and Arg407 in TM11; thus linking

the ICH domain to a TM involved in the inter-bundle salt-bridge network. In inward-facing GLUT5, these interactions are broken (Extended Data Fig. 7c) and, as observed in the inward-facing hGLUT1 and Xyle structures^{21,23,24}, ICH5 could not be built (Extended Data Fig. 2a). Therefore, the role of the ICH domain might be to provide additional stabilization of the outward-facing conformation, as suggested previously²¹.

TM7 and TM10 form substrate-induced gates

In GLUT5, the N- and C-terminal TM bundles undergo a small rotation of $\sim 15^\circ$ between the open outward- and inward-facing conformations (Fig. 1). As observed in other MFS transporter structures¹⁶, cavity-closing contacts in GLUT5 form mostly between TM1 and TM7 on the outside, and between TM4 and TM10 on the inside (Figs 1c and 4a). Among these helices, however, TM7 and TM10 in the C-terminal TM bundle seem to have the most prominent roles in occluding the substrate-binding site from the outside and inside, respectively (Fig. 4b). Comparisons of the substrate-free open GLUT5 structures with the substrate-occluded Xyle structures^{22,23} highlight the central role of TM7 and TM10 in gating (Extended Data Fig. 8). Inverted-symmetry-related TM7 and TM10 make up highly conserved sugar transporter signatures², and they also form a large fraction of the substrate-binding site in hGLUT1 (ref. 21) and Xyle (refs 21, 22), in agreement with previous functional data^{28,32,33}. Between the outward-open GLUT5 conformation and outward-occluded Xyle conformation²², the extracellular half of TM7 shows the largest shift towards the substrate-binding site (Extended Data Fig. 8a). Similarly, from the inward-occluded Xyle conformation²³ to the inward-open GLUT5 conformation, the intracellular half of TM10 shows the largest shift away from the substrate-binding site (Extended Data Fig. 8b, c); this TM10 movement has also been described previously for Xyle^{23,24}. In GLUT5 the observed conformational changes in TM7 and TM10 occur at hinge points that contain glycine residues (Fig. 4b, c and Extended Data Fig. 4). Tyr382 in TM10 interacts tightly with TM7 residues Ile295 and Val292 in the outward-facing state and these interactions do not take place in the inward-facing conformation (Fig. 4c). The interactions between TM7 and TM10 seem important to coordinate their conformational changes during transport, because mutation of the Ile295 equivalent residue in hGLUT5 to valine or alanine abolishes D-fructose transport³⁴. Consistently, similar interface-perturbing mutations of Ile295 significantly reduce D-fructose binding in rGLUT5 (Extended Data Fig. 6d). Furthermore, in hGLUT1 the equivalent residue to Val292 (Ile287) was substituted to every other amino acid, and each mutant shows markedly different D-glucose affinities and transport kinetics³⁵. Therefore, it suggests that the interactions between TM7 and TM10 are important and tuned with respect to substrate affinity and transport kinetics.

Conclusions

Symmetrical substrate binding and rigid-body movements of the N- and C-terminal TM bundles around a centrally located substrate-binding site form the structural basis of the ‘rocker-switch’ mechanism in MFS transporters^{16,36}. However, here we have described asymmetrical rearrangements in the MFS subfamily of sugar porters, consistent with the asymmetrical binding mode of the sugar reported previously in Xyle and hGLUT1 structures^{21,22}. Therefore, we conclude that transport in GLUT5 is not only controlled by the rocker-switch-type movement²¹ of the N- and C-terminal TM bundles, but also by a gated-pore mechanism, in the form of local, gating movements by TM7 and TM10 in the C-terminal TM bundle that are coupled to substrate binding and release (Fig. 5 and Supplementary Videos 1, 2 and 3). TM10 is also part of the inter-TM bundle salt-bridge network, indicating how local gating and global rocker-switch-type movements are coupled. Importantly, a deeper understanding of GLUT5 structure and

transport mechanism, as described here, should facilitate the structure-based design of novel inhibitors with therapeutic potential.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 19 February; accepted 14 July 2015.

Published online 30 September 2015.

- Mueckler, M. & Thorens, B. The SLC2 (GLUT) family of membrane transporters. *Mol. Aspects Med.* **34**, 121–138 (2013).
- Zhao, F. Q. & Keating, A. F. Functional properties and genomics of glucose transporters. *Curr. Genomics* **8**, 113–128 (2007).
- Simpson, I. A., Vannucci, S. J. & Maher, F. Glucose transporters in mammalian brain. *Biochem. Soc. Trans.* **22**, 671–675 (1994).
- James, D. E., Strube, M. & Mueckler, M. Molecular cloning and characterization of an insulin-regulatable glucose transporter. *Nature* **338**, 83–87 (1989).
- Burant, C. F., Takeda, J., Brot-Laroche, E., Bell, G. I. & Davidson, N. O. Fructose transporter in human spermatozoa and small intestine is GLUT5. *J. Biol. Chem.* **267**, 14523–14526 (1992).
- Kayano, T. et al. Human facilitative glucose transporters. Isolation, functional characterization, and gene localization of cDNAs encoding an isoform (GLUT5) expressed in small intestine, kidney, muscle, and adipose tissue and an unusual glucose transporter pseudogene-like sequence (GLUT6). *J. Biol. Chem.* **265**, 13276–13282 (1990).
- Blakemore, S. J. et al. The GLUT5 hexose transporter is also localized to the basolateral membrane of the human jejunum. *Biochem. J.* **309**, 7–12 (1995).
- Douard, V. & Ferraris, R. P. Regulation of the fructose transporter GLUT5 in health and disease. *Am. J. Physiol. Endocrinol. Metab.* **295**, E227–E237 (2008).
- Rand, E. B., Depaoli, A. M., Davidson, N. O., Bell, G. I. & Burant, C. F. Sequence, tissue distribution, and functional characterization of the rat fructose transporter GLUT5. *Am. J. Physiol.* **264**, G1169–G1176 (1993).
- Shepherd, P. R., Gibbs, E. M., Wesslau, C., Gould, G. W. & Kahn, B. B. Human small intestine facilitative fructose/glucose transporter (GLUT5) is also present in insulin-responsive tissues and brain. Investigation of biochemical characteristics and translocation. *Diabetes* **41**, 1360–1365 (1992).
- Zisman, A. et al. Targeted disruption of the glucose transporter 4 selectively in muscle causes insulin resistance and glucose intolerance. *Nature Med.* **6**, 924–928 (2000).
- Barone, S. et al. Slc2a5 (Glut5) is essential for the absorption of fructose in the intestine and generation of fructose-induced hypertension. *J. Biol. Chem.* **284**, 5056–5066 (2009).
- Douard, V. & Ferraris, R. P. The role of fructose transporters in diseases linked to excessive fructose intake. *J. Physiol.* **591**, 401–414 (2013).
- Zamora-Leon, S. P. et al. Expression of the fructose transporter GLUT5 in human breast cancer. *Proc. Natl Acad. Sci. USA* **93**, 1847–1852 (1996).
- Warburg, O. On respiratory impairment in cancer cells. *Science* **124**, 269–270 (1956).
- Yan, N. Structural advances for the major facilitator superfamily (MFS) transporters. *Trends Biochem. Sci.* **38**, 151–159 (2013).
- Madej, M. G., Sun, L., Yan, N. & Kaback, H. R. Functional architecture of MFS D-glucose transporters. *Proc. Natl Acad. Sci. USA* **111**, E719–E727 (2014).
- Maiden, M. C., Davis, E. O., Baldwin, S. A., Moore, D. C. & Henderson, P. J. Mammalian and bacterial sugar transport proteins are homologous. *Nature* **325**, 641–643 (1987).
- Pao, S. S., Paulsen, I. T. & Saier, M. H. Jr. Major facilitator superfamily. *Microbiol. Mol. Biol. Rev.* **62**, 1–34 (1998).
- Farwick, A., Bruder, S., Schadoweg, V., Oreb, M. & Boles, E. Engineering of yeast hexose transporters to transport D-xylose without inhibition by D-glucose. *Proc. Natl Acad. Sci. USA* **111**, 5159–5164 (2014).
- Deng, D. et al. Crystal structure of the human glucose transporter GLUT1. *Nature* **510**, 121–125 (2014).
- Sun, L. et al. Crystal structure of a bacterial homologue of glucose transporters GLUT1–4. *Nature* **490**, 361–366 (2012).
- Quistgaard, E. M., Low, C., Moberg, P., Tresaugues, L. & Nordlund, P. Structural basis for substrate transport in the GLUT-homology family of monosaccharide transporters. *Nature Struct. Mol. Biol.* **20**, 766–768 (2013).
- Wisdechairs, G., Park, M. S., Iadanza, M. G., Zheng, H. & Gonen, T. Proton-coupled sugar transport in the prototypical major facilitator superfamily protein XylE. *Nat. Commun.* **5**, 4521 (2014).
- Iancu, C. V., Zamoon, J., Woo, S. B., Aleshin, A. & Choe, J. Y. Crystal structure of a glucose/H⁺ symporter and its mechanism of action. *Proc. Natl Acad. Sci. USA* **110**, 17862–17867 (2013).
- Garcia, J. C., Strube, M., Leingang, K., Keller, K. & Mueckler, M. M. Amino acid substitutions at tryptophan 388 and tryptophan 412 of the HepG2 (Glut1) glucose transporter inhibit transport activity and targeting to the plasma membrane in *Xenopus* oocytes. *J. Biol. Chem.* **267**, 7770–7776 (1992).
- Li, Q. et al. Cloning and functional characterization of the human GLUT7 isoform SLC2A7 from the small intestine. *Am. J. Physiol. Gastrointest. Liver Physiol.* **287**, G236–G242 (2004).
- Mueckler, M. & Makepeace, C. Analysis of transmembrane segment 10 of the Glut1 glucose transporter by cysteine-scanning mutagenesis and substituted cysteine accessibility. *J. Biol. Chem.* **277**, 3498–3503 (2002).
- Andersson, M. et al. Proton-coupled dynamics in lactose permease. *Structure* **20**, 1893–1904 (2012).
- Law, C. J. et al. Salt-bridge dynamics control substrate-induced conformational change in the membrane transporter GlpT. *J. Mol. Biol.* **378**, 828–839 (2008).
- Schürmann, A. et al. Role of conserved arginine and glutamate residues on the cytosolic surface of glucose transporters for transporter function. *Biochemistry* **36**, 12897–12902 (1997).
- Seatter, M. J., De la Rue, S. A., Porter, L. M. & Gould, G. W. QLS motif in transmembrane helix VII of the glucose transporter family interacts with the C-1 position of D-glucose and is involved in substrate selection at the exofacial binding site. *Biochemistry* **37**, 1322–1326 (1998).
- Hruz, P. W. & Mueckler, M. M. Cysteine-scanning mutagenesis of transmembrane segment 7 of the GLUT1 glucose transporter. *J. Biol. Chem.* **274**, 36176–36180 (1999).
- Manolescu, A., Salas-Burgos, A. M., Fischbarg, J. & Cheeseman, C. I. Identification of a hydrophobic residue as a key determinant of fructose transport by the facilitative hexose transporter SLC2A7 (GLUT7). *J. Biol. Chem.* **280**, 42978–42983 (2005).
- Kasahara, T., Maeda, M., Boles, E. & Kasahara, M. Identification of a key residue determining substrate affinity in the human glucose transporter GLUT1. *Biochim. Biophys. Acta* **1788**, 1051–1055 (2009).
- Karpowich, N. K. & Wang, D. N. Structural biology. Symmetric transporters for asymmetric transport. *Science* **321**, 781–782 (2008).
- Radestock, S. & Forrest, L. R. The alternating-access mechanism of MFS transporters arises from inverted-topology repeats. *J. Mol. Biol.* **407**, 698–715 (2011).
- Solcan, N. et al. Alternating access mechanism in the POT family of oligopeptide transporters. *EMBO J.* **31**, 3411–3421 (2012).
- Fukuda, M. et al. Structural basis for dynamic mechanism of nitrate/nitrite antiport by NarK. *Nature Commun.* **6**, 7097 (2015).

Supplementary Information is available in the online version of the paper.

Acknowledgements We are grateful to D. Slotboom, A. Cameron and S. Newstead for discussions and comments, and J. Mansfield for assistance with large-scale yeast fermentations and H. Unno with rGLUT5 crystallization. Data were collected at the European Synchrotron Radiation Facility, Diamond Light Source, and SPring-8 (proposal numbers 2011A1393, 2011B1229, 2012A1184, 2012B1253, 2013A1241, 2013B1237, 2014A1348 and 2014B1407), with assistance from beamline scientists. This work was funded by the Knut and Alice Wallenberg Foundation (D.D.), The Royal Society through the University Research Fellow scheme (D.D.), the BBSRC (BB/G02325/1 to S.I.), the ERATO Human Receptor Crystallography Project of the Japan Science and Technology Agency (JST) (S.I.), by the Research Acceleration Program of the JST (S.I.), by the Targeted Proteins Research Program of the Ministry of Education, Culture, Sports, Science and Technology (MEXT) of Japan (S.I.), and by Grants-in-Aids for Scientific Research from the MEXT (No. 22570114 to N.N.), and by the Platform for Drug Discovery, Informatics, and Structural Life Science from the MEXT (T.K.). The authors are grateful for the use of the Membrane Protein Laboratory funded by the Wellcome Trust (grant 062164/Z/00/Z) at the Diamond Light Source Limited, and The Centre for Biomembrane Research (CBR), supported by the Swedish Foundation for Strategic Research. H.J.K. was a recipient of a Human Frontiers Postdoctoral fellowship and D.D. acknowledges support from EMBO through the Young Investigator Program (YIP).

Author Contributions N.N., S.I. and D.D. designed the project. Cloning, expression screening and initial crystallization of rat and bovine GLUT5 was carried out by H.J.K., Y.S. and D.D. Crystal optimization of bovine GLUT5 was carried out by H.J.K. and G.V. Data collection, structure determination and refinement of bovine GLUT5 was carried out by G.V. Generation of rat GLUT5 scFv fragment was carried out by N.N., Y.N., T.M., Y.N.-N., O.K.-A., H.I., T.A., T.K. and T.Ha. Expression and purification of the Fv fragment was carried out by N.N., Y.N., Y.Sa., H.A. and T.Hi. Co-crystallization of rat GLUT5–Fv complex and data collection was performed by N.N. and Y.N. with assistance from T.Hi. and S.I. Structure determination and refinement of rat GLUT5–Fv was carried out by T.S. Experiments for functional analysis were designed by M.K. and D.D. and carried out by M.K., D.D., S.A.H. and A.A.Q. Modelling of GLUT5 was carried out by M.C. The manuscript was prepared by N.N., H.J.K., G.V., S.I. and D.D. All authors discussed the results and commented on the manuscript.

Author Information The coordinates and the structure factors for bovine and rat GLUT5 have been deposited in the Protein Data Bank under accessions 4YB9 and 4YBQ, respectively. Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to N.N. (nnomura@mfour.med.kyoto-u.ac.jp), S.I. (siwata@mfour.med.kyoto-u.ac.jp) or D.D. (ddrew@dbb.se).

METHODS

No statistical methods were used to predetermine sample size.

Rattus norvegicus GLUT5 cloned full-length sequence (UniProt accession number P43427); N50Y deglycosylation mutation is underlined and additional C-terminal residues retained after TEV cleavage are shown in italic (see next section for cloning details): MEKEDQEKTKGLTLVLALATFLAAGSSSQ YGYNVA AVNSPSEFMQQFYDYTYDRNKENIESFTLTLLWSLTVMFPFG GFIGSLMVGFLVNNLGRKGALLFNNIFSILPAILMGCSKIAKSEIIHARSLLV GICAGISSNVVPMYLGELAPKNLRGALGVVPQLFITV GILVAQLFGLRSVL ASEEGWPILLGLTGVPAGLQLLLPFFPESPRYLLIQKKNESA AEKALQTLR GWKDVDMEMEEIRKEDEAEKAGFISVWKLFRMQSLRWQLISTIVLMT GQQLSGVNAIYYYADQIYLSAGVKSNDVQYV TAGTGAVNVFMTMVTV FVVELWGRNRLLIGFSTCLTACIVLTV ALALQNTISWMPYVIVCVIVY VIGHAVGPSPIALFITEIFLQSSRPSAYMIGGSVHWLSNFIVGLIFFIQQVG LGPYSFIHAIICLLTSIYIMVVPETKGRTFVEINQIFAKKNKVSDDVYPEKE EKELNDLPATREQENLYFQ.

Bos taurus GLUT5 cloned sequence containing 1–473 out of 501 residues (UniProt accession number P58353); N51A deglycosylation mutation is underlined and additional C-terminal residues retained after TEV cleavage are shown in italic (see next section for cloning details): MEPQDPVKREGRLTPVIVLAT LIAAFGSSSQYGYNVA AVNSPSEFMKDFYAYTYTYDRVGEYMNEFYLTLLW SVTVSMFPFGGLGSLMVGPLVNNLGRKGTLFNNIFSIVPALLMGFSELA KSFEMIIVARVLV GICAGLSSNVVPMYLGELAPKNWRGALGVVPQLFITI GILVAQIFGLRSLLANEEGWPIILLGTGIPAVLQLLFPFFPESPRYLLIQKK DEAAKASALRRRLRGWHDVDAIEIELEEDRAEKAVGFISVLKFKMRSLR WQVISIIVLMAGQSLSGVNAIYYADQIYLSAGVNEEDDVQYV TAGTGAV NVLITVCAIFVVELMGRRLLLGLFSVCFTACCVLTGALALQDVISWMPY VSIACVISYVIGHALGPSPIALLVTEIFLQSSRPAAYMVAGTVHWSNFTV GLVFFFIQVGLGAYSFVIFAVICLLTTVYIFLIIPETKSKTFIEINRENLYFQ.

Fv light chain variable region (V_L); additional C-terminal residues, retained after TEV cleavage, are shown in italics: ELDIVLTQSPSLPVLGDAQSIS RSSQSVHSNGNTYLEWYLQKPGQSPKLLIYKVSNRFGVDPDRFSGSGSGTD FTLKISRVEAEDLVGYCYCFQGSHPVPTFGGGTKLEIKTSENLYFQ.

Fv heavy chain variable region (V_H); additional C-terminal residues retained after TEV cleavage are shown in italics: LEVNLVESGGGLVQPGGSRKLS CAASGFTFSFGMHVWRQAPKGLVWVAHSSGSRITIDYADTVKGRFTIS RDNPKNTLFLQMTSLRSEDTAIYCCARGNGYDADLYWGQGSVTVSSA KTTTPSVTSENLYFQ.

Target identification using fluorescence-based screening methods. GLUT5 homologues were cloned into the GAL1 inducible TEV cleavable GFP–His₈ 2 μ vector pDDGFP2, transformed into the *Saccharomyces cerevisiae* strain FGY217 (MATa, ura3–52, lys2 Δ 201, and pep4 Δ)⁴⁰ and overexpressed as described previously^{41,42}. In brief, 10 ml *S. cerevisiae* FGY217 cell cultures in –URA media and 0.1% glucose were grown at 30 °C and expression was induced by the addition of final 2% (v/v) D-galactose at an OD₆₀₀ of 0.6. After ~22 h, cells were harvested, resuspended in 1 \times PBS buffer and overexpression levels assessed by whole-cell GFP fluorescence^{41,42}. Fusions with detectable expression levels were re-grown in larger 2-l culture volumes and membranes subsequently isolated. The monodispersity of expressed fusions was screened in crude dodecyl- β -D-maltopyranoside (DDM), decyl- β -D-maltopyranoside (DM), nonyl- β -D-maltopyranoside (NM) and *n*-dodecyl-*N,N*-dimethylamine-*N*-oxide (LDAO) solubilized membranes by fluorescence-detection size exclusion chromatography (FSEC)⁴³ as described previously⁴¹. Out of the GLUT5 homologues screened, rat and bovine GLUT5 showed the sharpest monodispersity profiles and was the most stable after purification in detergent⁴⁴.

Large-scale production and purification of rat and bovine GLUT5. For rat GLUT5, cells were harvested from 10 l *S. cerevisiae* cultures, resuspended in buffer containing 50 mM Tris-HCl pH 7.6, 1 mM EDTA, 0.6 M sorbitol, and lysed by mechanical disruption as described previously⁴². Membranes were isolated by ultracentrifugation at 195,000 g for 3 h, homogenized in 20 mM Tris-HCl pH 7.5, 0.3 M sucrose, 0.1 mM CaCl₂, frozen in liquid nitrogen and stored at –80 °C. Rat GLUT5 membranes were solubilized with 1% DDM in equilibration buffer (EB) consisting of 1 \times PBS, 10 mM imidazole, 150 mM NaCl, 10% glycerol. After ultracentrifugation at 195,000 g for 45 min, the supernatant was incubated with 10 ml of Ni²⁺-nitrilotriacetate affinity resin (Ni-NTA; Qiagen) for 2 h at 4 °C. The resin was washed with 100 ml of EB containing 0.05% DDM and 30 mM imidazole, and the protein eluted in 25 ml of EB containing 0.05% DDM and 250 mM imidazole. The eluted protein was incubated with TEV–His₈ protease to remove the GFP–His₈ tag and dialysed against 20 mM Tris-HCl pH 7.5, 150 mM NaCl, 10% glycerol, 0.03% DDM. Dialysed sample was then loaded onto a 5 ml HisTrap column (GE Healthcare) equilibrated in dialysis buffer, and the flow-through was collected and concentrated. The protein was further purified by size-exclusion chromatography

(SEC) using a Superdex 200 column (GE Healthcare) in buffer consisting of 20 mM Tris-HCl pH 7.5, 150 mM NaCl, 0.02% DDM. All rat GLUT5 mutants were generated with a standard PCR-based strategy and were sub-cloned, overexpressed and purified as described for wild type.

Protein complexes were prepared by incubation of the HisTrap-purified rat GLUT5 with the SEC-purified Fv fragment (supplemented with 0.02% (w/v) DDM) at a molar ratio of 1:1.5 for 1 h on ice. The complex was subjected to SEC (Superdex 200, GE Healthcare). The SEC step was repeated twice to ensure reproducibility of crystals. Peak fractions containing rat GLUT5–Fv complex were concentrated to ~10 mg ml^{–1} by ultrafiltration (Millipore, MWCO 50 kDa), and immediately used for crystallization experiments.

For bovine GLUT5, membranes were isolated from 40 l *S. cerevisiae* cultures. After membrane solubilization and ultracentrifugation (as described for rat GLUT5), the supernatant was supplemented with 55 mM imidazole, and incubated with 30 ml of Ni-NTA resin for 2 h at 4 °C. The resin was washed with 600 ml of EB supplemented with 55 mM imidazole and 0.03% DDM, and the protein was eluted using 125 ml of EB containing 250 mM imidazole and 0.03% DDM. After cleavage of GFP–His₈ cleavage by TEV–His₈ protease the material was re-loaded onto the same 30 ml Ni-NTA resin column (no imidazole), the flow-through was concentrated and passed through a 1 ml Ni-NTA column for further clean up. The untagged protein was further purified by SEC in buffer containing 10 mM Tris-HCl, pH 7.5, 150 mM NaCl, and 0.09% undecyl- β -maltopyranoside (UDM β), and then concentrated to ~4 mg ml^{–1} for crystallization experiments.

Generation of single-chain Fv (scFv) fragments. Animal experiments conformed to the guidelines outlined in the Guide for the Care and Use of Laboratory Animals of Japan and were approved by the University of Tokyo Animal Care Committee (approval no. RAC07101). Rat GLUT5-specific scFv fragment were generated essentially as previously described⁴⁵. In brief, purified GLUT5 was reconstituted into liposomes containing chicken egg yolk phosphatidylcholine (egg PC; Avanti) and adjuvant lipid A (Sigma-Aldrich) by detergent removal method. Small unilamellar proteoliposomes were prepared by sonication. MRL/lpr mice were immunized with 0.1 mg of the proteoliposome antigen three times at two-week intervals. Immunized mice were killed, and RNA in their splenocytes was isolated and converted into cDNA via reverse-transcription PCR. The V_L and V_H repertoire was assembled via an 18-amino-acid flexible linker and cloned into the phage-display vector pComb3XSS. Biotinylated proteoliposomes were prepared by reconstituting GLUT5 with a mixture of egg PC and 1,2-dipalmitoyl-sn-glycero-3-phosphoethanolamine-*N*-(cap biotinyl) (16:0 biotinyl Cap-PE; Avanti), and used as binding targets for scFv-phage selection. The targets were immobilized onto streptavidin-coated paramagnetic beads (Dynabeads) or streptavidin-coated microplates (Nunc). After four rounds of biopanning, proteoliposome-targeted enzyme-linked immunosorbent assays ('liposome ELISAs') were performed on periplasmic extracts of individual colonies. Positive clones were collected and evaluated using a Biacore T100 (GE Healthcare).

Expression and purification of Fv fragment. Antibody fragments in the scFv format are undesirable for use as crystallization chaperones because they are able to intermolecularly form domain-swapped dimers, and dimer-monomer equilibrium may increase structural heterogeneity. Therefore, we used Fv fragments for crystallization trials. The Fv fragments were expressed in *Brevibacillus choshinensis* as a secreted His₈-tagged protein and purified from culture medium. The cells were grown at 30 °C with 200 rpm in 2SY medium (soypton 40 g l^{–1}, yeast extract 5 g l^{–1}, glucose 20 g l^{–1}, CaCl₂ 0.15 g l^{–1}) supplemented with 50 mg l^{–1} neomycin. The cells expressing the V_L and V_H were initially grown separately as overnight pre-cultures. The pre-cultures were then combined and diluted in 2SY medium to give OD₆₀₀ of 0.02 of each strain and grown for further 65–70 h. The cells were removed by centrifugation at 6,000 g for 15 min. The recovered culture supernatant was adjusted to a final ammonium sulfate concentration of 60% saturation, and the precipitate was pelleted, dissolved in TBS buffer (10 mM Tris-HCl, pH 7.5, 150 mM NaCl), and dialysed overnight against the same buffer. Dialysed sample was mixed with Ni-NTA resin equilibrated with Fv1 buffer (10 mM Tris-HCl, pH 7.5, 150 mM NaCl, 20 mM imidazole). Bound proteins were eluted with Fv2 buffer (10 mM Tris-HCl, pH 7.5, 150 mM NaCl, 250 mM imidazole), mixed with TEV–His₈ and dialysed overnight against TBS buffer. Cleaved His₈ tag and TEV–His₈ were removed by binding to a HisTrap column equilibrated with Fv1 buffer. The tag-free Fv fragment was concentrated and loaded onto a HiLoad16/60 Superdex 75 column (GE Healthcare) equilibrated with TBS buffer. Peak fractions were pooled, concentrated, flash frozen in liquid nitrogen, and stored at –80 °C.

Transport activity of reconstituted GLUT5. The proteoliposome D-fructose uptake assay was modified from that previously described for human GLUT1 (ref. 46). In brief, purified rat GLUT5 was reconstituted by the freeze-thaw/extrusion method. Crude lipids were extracted from bovine liver and sonicated to make unilamellar liposomes. 500 μ l of a mixture containing ~10 μ g of purified GLUT5 and 20 mg of liposomes in 10 mM TrisSO₄ (pH 7.5) was flash frozen and

thawed at room temperature. Large, unilamellar proteoliposomes were prepared by extrusion (LiposoFast, Avestin; membrane pore size, 400 nm). For each time point, 10 μ l proteoliposomes (0.4 mg lipid; 0.2 μ g GLUT5) was added to 10 μ l transport buffer containing 10 mM TrisSO₄ and 2 mM MgSO₄ (pH 7.5) with or without the addition of 0.5 mM HgCl₂. Time course of 0.1 mM [¹⁴C]-D-fructose transport was measured at 25 °C at the indicated time intervals and stopped by the addition of cold buffer containing 10 mM TrisSO₄, 2 mM MgSO₄ (pH 7.5) and 0.5 mM HgCl₂ and immediately filtered. Non-specific uptake was estimated with 0.1 mM [¹⁴C]-L-glucose. The radioactivity corresponding to the internalized substrate was measured by scintillation counting. Each experiment was performed in triplicate and data points shown indicate average values of two technical replicates.

Substrate specificity. Unless otherwise stated the rat GLUT5 deglycosylation mutant (N50Y) used for structure determination is referred to as wild-type-like protein (WT). Only rat GLUT5 wild type and mutants showing a monodisperse peak in DDM during gel filtration were assessed for their ability to bind sugar by tryptophan fluorescence. In all experiments, either purified rat GLUT5 wild type or mutants were diluted to 0.06 mg ml⁻¹ in purification buffer containing 150 mM NaCl₂, 20 mM Tris pH 7.5, 0.03% DDM. D-fructose, L-fructose, D-glucose, D-mannose, D-xylose and D-galactose stocks, that were freshly prepared in purification buffer, were added to the diluted protein to reach a final concentration of 40 mM. After each sugar addition the sample was incubated at room temperature for 2 min before measuring tryptophan fluorescence. Measurements were performed using a Cary Eclipse Fluorescence Spectrophotometer (Agilent Technologies), with an excitation wavelength of 295 nm and the emission spectra recorded in the range of 300–400 nm with a 5 nm excitation slit in a High Precision Cell Quartz Suprasil 10 × 2 mm cuvette (Hellma Analytics). The emission peak of 338 nm was taken as an average of the four wavelengths 337.03, 337.96, 339.06 and 340 nm. Each experiment was carried out in triplicate in the absence and presence of the well-known GLUT inhibitor HgCl₂ (5 mM)⁴⁷. To compare mutant binding to D-fructose, nonspecific tryptophan fluorescence quenching of 2.5% (as measured with 40 mM L-fructose) was first subtracted from both wild-type and mutant quenching levels before calculating the final percentage of quenching relative to wild type.

Binding analysis. D-fructose or D-glucose stocks, that were freshly prepared in purification buffer, were sequentially added to the diluted protein to a final concentration of 40 mM. After each sugar addition the sample was incubated at room temperature for 2 min before measuring. Tryptophan fluorescence measurements were performed as above. Each experiment was carried out in triplicate in the absence and presence of the known GLUT inhibitor HgCl₂ (5 mM). The rat GLUT5 curves were fitted by nonlinear regression using the software Prism. The bovine GLUT5 construct used for structural studies binds D-fructose with similar affinity as rat GLUT5 with a K_d of 5.5 ± 0.6 mM.

Crystallization. The structure of the Fv fragment was used as a search model for the molecular replacement. Crystals of the Fv fragment appeared in the well buffer at 20 °C containing 0.1 M CAPS-NaOH (pH 10.5), 1.6 M (NH₄)₂SO₄, 0.1 M Li₂SO₄, 6.5 mM n-nonyl- β -D-glucoside.

For rat GLUT5, crystals of rGLUT5–Fv complex (in DDM) used for structural determination were grown at 20 °C by hanging-drop vapour diffusion. A 400 μ l reservoir containing 33–35% PEG 400, 0.12 M CaCl₂, 0.1 M Tris–HCl (pH 8.0) was equilibrated against a 2 μ l drop containing a 1:1 mixture of the complex and reservoir solution. After 2–3 weeks of growth, crystals were dehydrated by step-wise equilibration of the drops against 400 μ l reservoirs containing increasing concentration of PEG 400, in steps of 5% and up to a final concentration of 70%. Crystals were then flash-frozen and stored in liquid nitrogen.

For bovine GLUT5, crystals were grown over the course of 2–3 weeks at 4 °C in 2 μ l drops (500 μ l well) consisting of the well solution 31–35% PEG 300, 0.125 M HEPES pH 8.0, 0.125 M NaCl, 0.125 M LiSO₄, mixed first with 70 mM HEGA-10 (4:1) and then the protein solution (1:1), and placed over the well solution diluted with water (4:1) before sealing. For freezing in liquid nitrogen, crystals were soaked for 1 min in the mother liquor supplemented with ~15% PEG 300.

Data collection, structure determination and analysis. For the Fv fragment, diffraction data were collected at 100 K at Spring-8 beamline BL41XU (Japan) and processed using HKL2000 packages⁴⁸ and the CCP4 suite⁴⁹. The crystal belonged to the space group C2 with two molecules per asymmetric unit. The structure was determined by molecular replacement with Molrep⁵⁰ using the Fv portion of an antibody structure (PDB code 1IGC) as a search model. Refinement was performed till the R_{free} value decreased to ~22% with REFMAC5. For rat GLUT5–Fv complex, diffraction data were collected at 100 K at Spring-8 beamline BL41XU (Japan) and processed using HKL2000 packages⁴⁸ and the CCP4 suite⁴⁹. The data set used for structure determination and refinement was generated by the combination of 4 data sets from 4 independent crystals. The space group was determined to be P2₁, with two rat GLUT5–Fv complexes per asym-

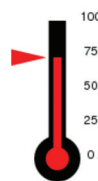
metric unit. The structure was determined by molecular replacement with Molrep⁵⁰ using two search models (polyalanine of the transmembrane region in the outward-facing conformation of FucP (PDB code 3O7Q), and the separately determined 1.5 Å structure of the Fv fragment used in this study). Refinement was performed with PHENIX⁵¹ followed by manual examination and rebuilding of the refined coordinates using COOT⁵². Recently determined structures of Xyle (PDB code 4GC0) and human GLUT1 (PDB code 4PYP) helped with the modelling. The 6 N-terminal residues (Met1–Gln6), 22 C-terminal residues (Ser481–Gln502), and 22 residues (Asn39–Asn60) in molecule A and 12 residues (Met45–Arg56) in molecule B in TM2 are not included in the structure as they did not have interpretable densities. For bovine GLUT5, data were collected on frozen crystals at beamline I02, Diamond Light Source (UK). The data set used for structure determination and refinement was generated using HKL2000 by processing and scaling together two data sets collected on two different parts of the same crystal, and by correcting for anisotropy (UCLA server, <http://services.mbi.ucla.edu/anisocore/>) (Extended Data Tables 1 and 2 and Extended Data Fig. 1). The structure was solved by molecular replacement using the N-terminal and C-terminal bundles of rat GLUT5 as independent search models in PHASER, and refined using REFMAC5 and BUSTER against data up to 3.0 Å resolution after anisotropy correction⁵³, with rounds of re-building in COOT⁵². Final refinement was performed using PHENIX⁵¹ at 3.2 Å resolution (Extended Data Fig. 1). Structural alignments were performed using the align command of PYMOL software using C α coordinates.

Homology modelling. Models for the outward-occluded and inward-occluded conformations were based on corresponding Xyle crystal structures (4GBY and 4JA3, respectively) and were produced using Modeller 9v12 (ref. 54). Multiple alignment of GLUT5 homologues (*Homo sapiens*, *Canis lupus*, *Rattus norvegicus*, *Felis catus*, *Gallus gallus*, *Anolis carolinensis*) and xyle sequences were combined with structural data from the Xyle structures and the GLUT5 structures using expresso mode of t-coffee⁵⁵. The multiple alignment in combination with helical restraints corresponding to the TMs of bovine GLUT5 structure was used as an input for Modeller. For each conformation, 20 models were generated with loop optimization. The DOPE scoring function was used to select the final model. Outward-facing bovine GLUT5 was modelled with the same protocol based on the rat outward-facing GLUT5 structure and alignment of GLUT5 sequences. The four conformations (outward-open, outward-occluded, inward-occluded and inward-open) were morphed and rendered using PyMOL (<http://www.pymol.org/>).

40. Kota, J., Gilstring, C. F. & Ljungdahl, P. O. Membrane chaperone Shr3 assists in folding amino acid permeases preventing precocious ERAD. *J. Cell Biol.* **176**, 617–628 (2007).
41. Newstead, S., Kim, H., von Heijne, G., Iwata, S. & Drew, D. High-throughput fluorescent-based optimization of eukaryotic membrane protein overexpression and purification in *Saccharomyces cerevisiae*. *Proc. Natl Acad. Sci. USA* **104**, 13936–13941 (2007).
42. Drew, D. et al. GFP-based optimization scheme for the overexpression and purification of eukaryotic membrane proteins in *Saccharomyces cerevisiae*. *Nature Protocols* **3**, 784–798 (2008).
43. Kawate, T. & Gouaux, E. Fluorescence-detection size-exclusion chromatography for precrystallization screening of integral membrane proteins. *Structure* **14**, 673–681 (2006).
44. Sonoda, Y. et al. Benchmarking membrane protein detergent stability for improving throughput of high-resolution X-ray structures. *Structure* **19**, 17–25 (2011).
45. Suharni, et al. Proteoliposome-based selection of a recombinant antibody fragment against the human M2 muscarinic acetylcholine receptor. *Monoclon. Antib. Immunodiagn. Immunother.* **33**, 378–385 (2014).
46. Sonoda, Y. et al. Tricks of the trade used to accelerate high-resolution structure determination of membrane proteins. *FEBS Lett.* **584**, 2539–2547 (2010).
47. Sarkar, H. K., Thorens, B., Lodish, H. F. & Kaback, H. R. Expression of the human erythrocyte glucose transporter in *Escherichia coli*. *Proc. Natl Acad. Sci. USA* **85**, 5463–5467 (1988).
48. Otwinowski, Z. & Minor, W. Processing of X-ray diffraction data collected in oscillation mode. *Methods Enzymol.* **276**, 307–326 (1997).
49. Collaborative Computational Project, Number 4. The CCP4 suite: programs for protein crystallography. *Acta Crystallogr. D* **50**, 760–763 (1994).
50. Vagin, A. & Teplyakov, A. Molecular replacement with MOLREP. *Acta Crystallogr. D* **66**, 22–25 (2010).
51. Adams, P. D. et al. PHENIX: building new software for automated crystallographic structure determination. *Acta Crystallogr. D* **58**, 1948–1954 (2002).
52. Emsley, P. & Cowtan, K. Coot: model-building tools for molecular graphics. *Acta Crystallogr. D* **60**, 2126–2132 (2004).
53. Blanc, E. et al. Refinement of severely incomplete structures with maximum likelihood in BUSTER-TNT. *Acta Crystallogr. D* **60**, 2210–2221 (2004).
54. Eswar, N. et al. Comparative protein structure modeling using MODELLER. *Curr. Protocols Protein Sci.* **Ch. 2**, Unit 2.9. (2007).
55. Notredame, C., Higgins, D. G. & Heringa, J. T-Coffee: A novel method for fast and accurate multiple sequence alignment. *J. Mol. Biol.* **302**, 205–217 (2000).

a

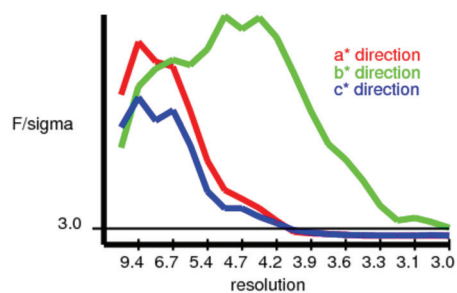
bGLUT5



Your data has SEVERE anisotropy based on the spread in values of the three principle components= 73.43 Å⁻²

The principle components are the exponential scale factors used to correct for anisotropy. They may be regarded as B factors applied to the three principle directions of the data set. Larger [values] indicate stronger anisotropy.

In decreasing order the 3 components are: 33.89 5.64 -39.53 Å⁻²



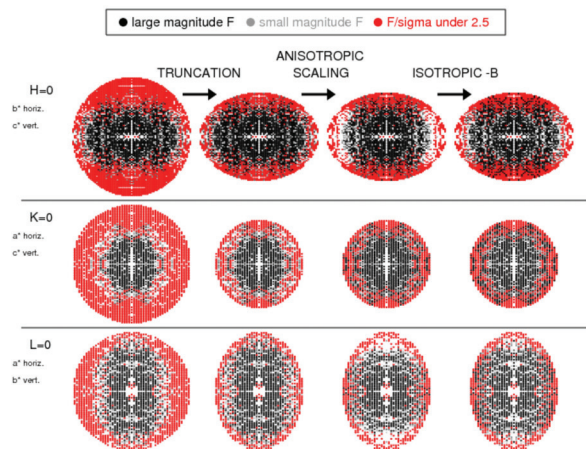
The recommended resolution limits along a*,b*,c* are

4.0 Å 3.0 Å 4.0 Å

These are the resolutions at which F/sigma drops below an arbitrary cutoff of 3.0

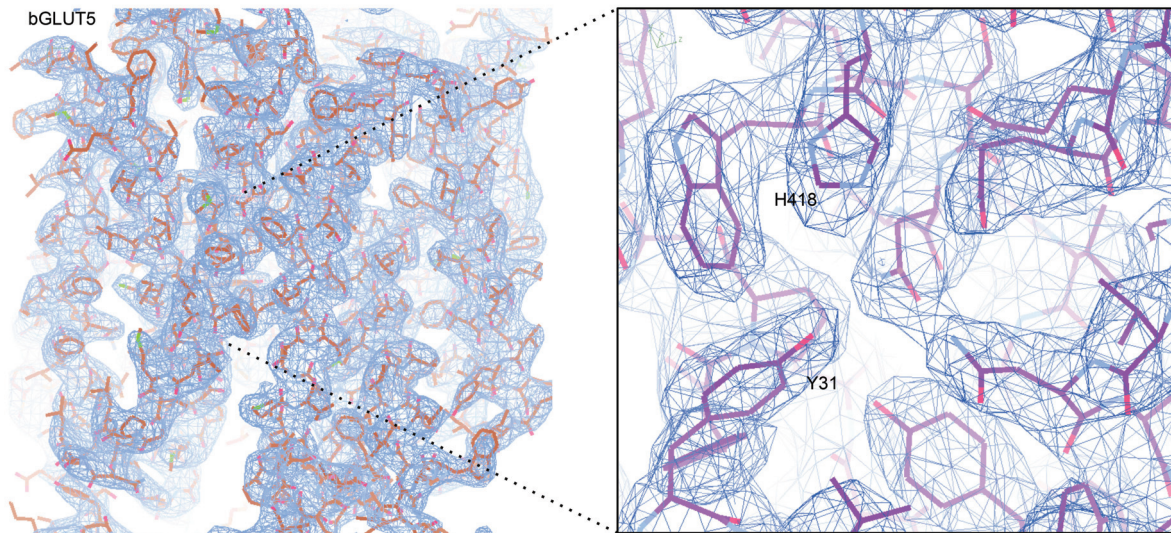
24188 reflections were in the initial data set. 10597 were discarded because they fell outside the specified ellipsoid with dimensions 1/4.0, 1/3.0, 1/4.0 Å along a*,b*,c*, respectively. These discarded reflections had an average F/sigma of 2.41.

13591 reflections remain after ellipsoidal truncation. Anisotropic scale factors were then applied to remove anisotropy from the data set. Lastly, an isotropic B of -81.00 Å⁻² was applied to restore the magnitude of the high resolution reflections diminished by anisotropic scaling. The following pseudo precession images illustrate the individual steps.



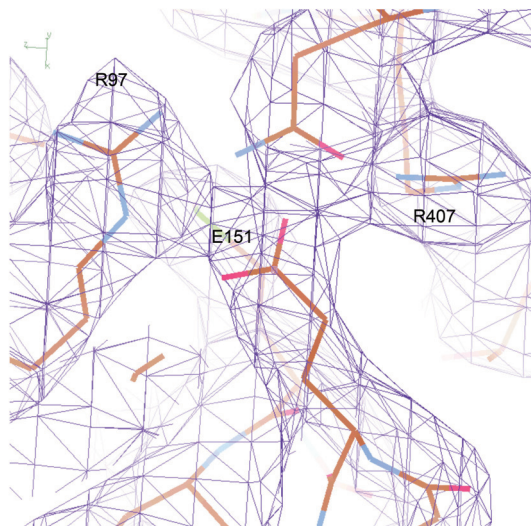
b

bGLUT5



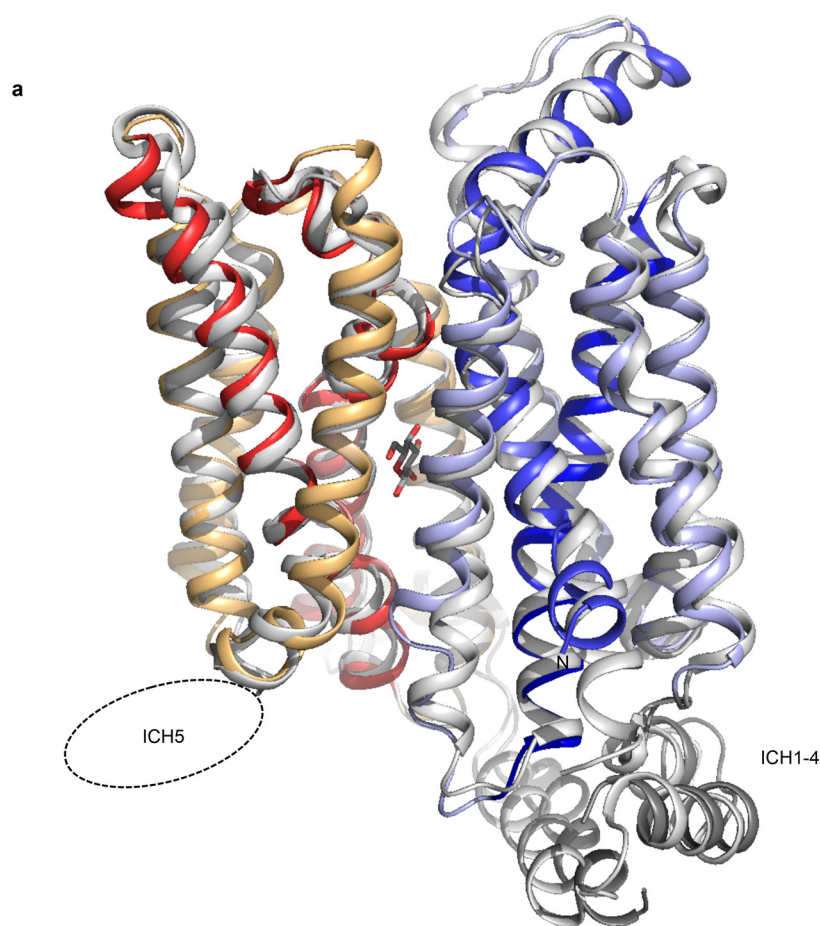
c

rGLUT5



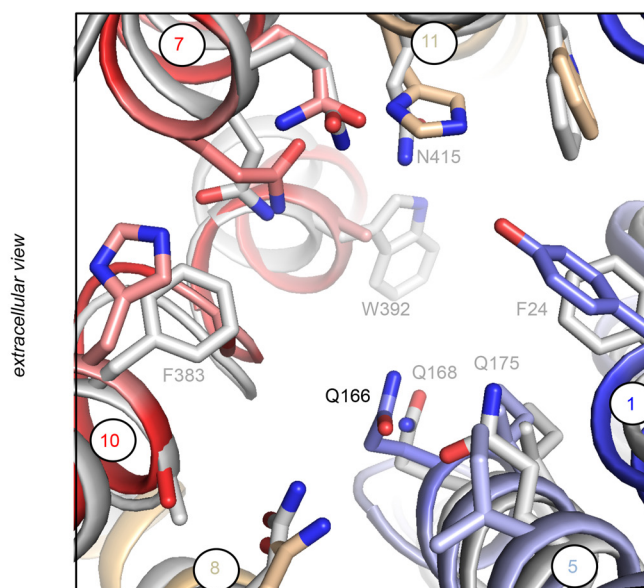
Extended Data Figure 1 | Anisotropy descriptors of bGLUT5 data reported by the UCLA-MBI Diffraction Anisotropy Server and $2F_o - F_c$ electron density maps for the bovine and rat GLUT5 structures. **a**, Degree of anisotropy of bGLUT5 data, resolution limits for the 3 principal axes (left), and panel illustrating steps along correction of bGLUT5 data for anisotropy (right).

b, Representative portions of the electron density map (1.5σ) for bGLUT5 overall model (left) and a close-up of the substrate binding site (right); residues highlighted are numbered based on rGLUT5 for the sake of clarity. **c**, Electron density (1.0σ) for rGLUT5 showing one of the inter-bundle salt-bridge clusters that form in the open outward-facing conformation.



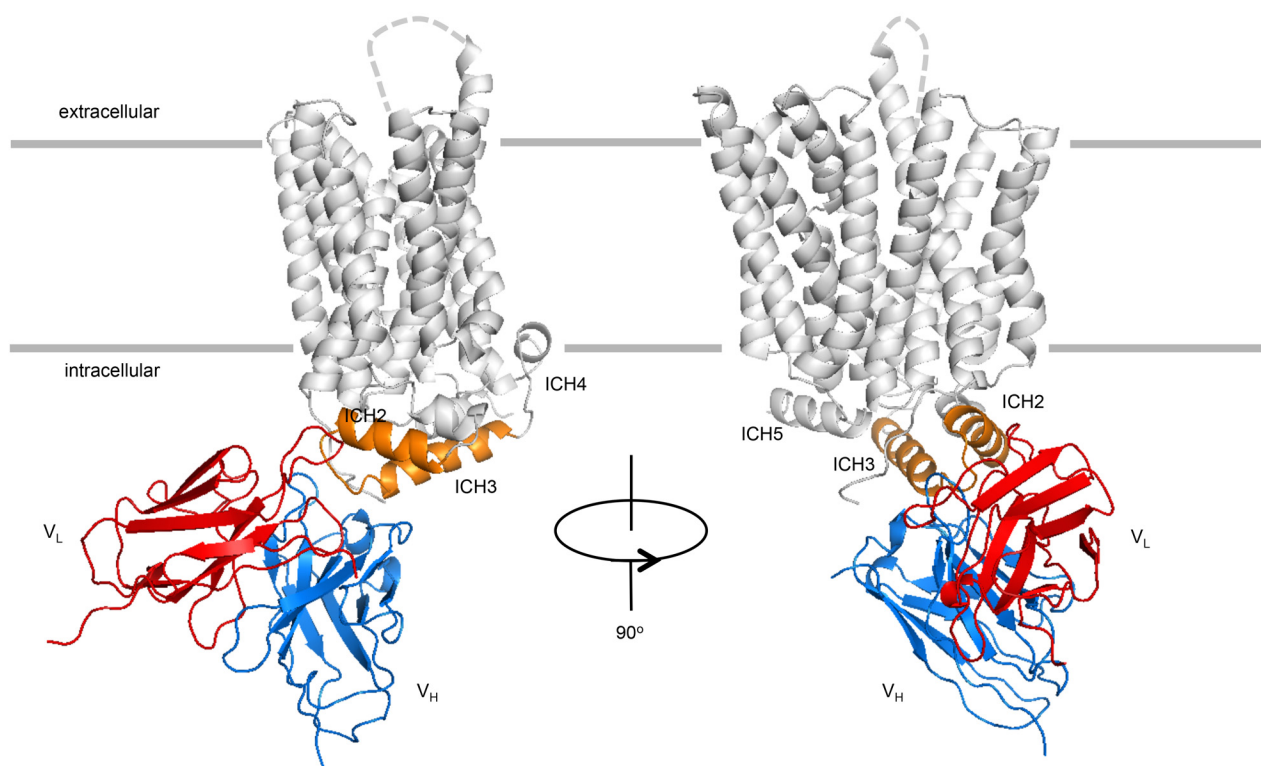
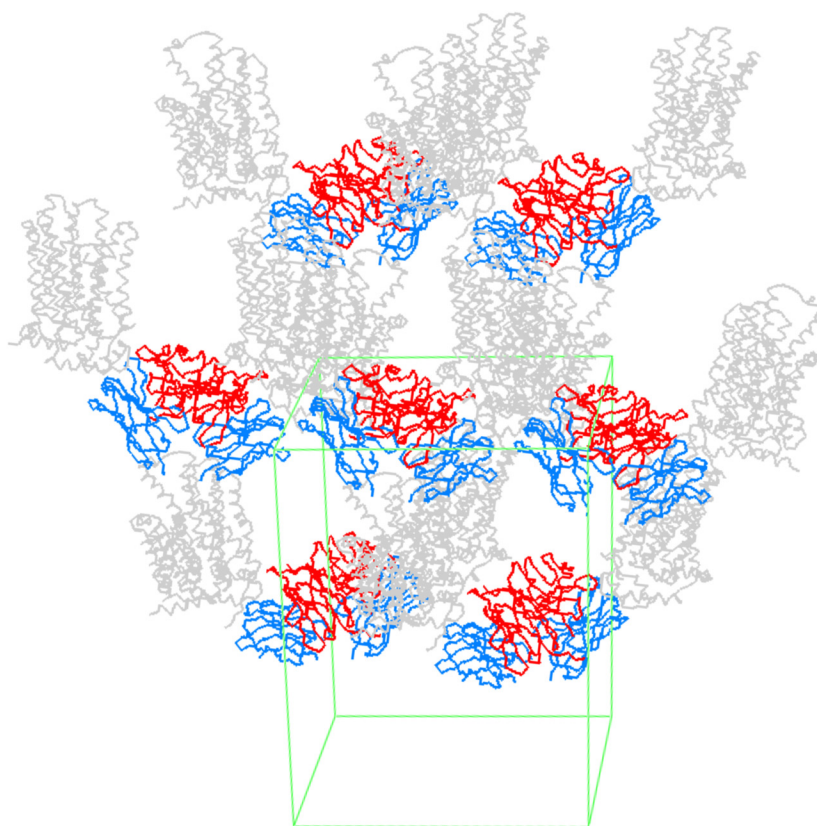
b

GLUT5 and Xyle (grey)



Extended Data Figure 2 | Superimposition of open inward-facing bGLUT5 and hGLUT1 structures, and comparison of the substrate-binding site in bGLUT5 and inward-facing Xyle **a**, Ribbon representation of inward-facing bGLUT5 (coloured as in Fig. 1a) and inward-facing hGLUT1 (light grey) structures, as viewed in the plane of the membrane. The D-glucopyranoside moiety of the detergent molecule bound to GLUT1 (*n*-nonyl- β -D-glucopyranoside (β -NG)) is shown in stick representation. Density for ICH5 at the C terminus is missing in both hGLUT1 and bGLUT5 inward-facing

structures and highlighted with the dotted ellipse. The beginning of TM1 kinks further outwards in the bGLUT5 structure compared to hGLUT1 and residues 1–18 could not be built. The r.m.s.d. (root mean square deviation) after superposition of the two structures is 1.12 Å for 364 pairs of C α atoms (see Methods). **b**, The substrate-binding in the inward-facing bGLUT5 structure (coloured as in Fig. 1) is very similar to that seen in inward-facing Xyle (4JA4) structure (shown in light grey). Only non-conserved residues and the equivalent glutamine to Q166 are labelled for Xyle.

a**b****Extended Data Figure 3 | Structure of the rat GLUT5-Fv complex.**

a, Cartoon representation of the complex between rGLUT5 (grey) and 4D111Fv (heavy-chain variable region (V_H) is in blue; light-chain variable region (V_L) is in red). 4D111Fv binds to the cytoplasmic domain of GLUT5,

including ICH2 (residues 226, 230, 234), the loop between ICH2 and ICH3 (residues 238, 240, 241), and ICH3 (residue 243), with $\sim 848 \text{ \AA}^2$ of buried surface area at the interface. **b**, Packing of the rat GLUT5-Fv complex molecules in the crystal. The unit cell is represented as green lines.



rGLUT5 1 MEKEDQ-----E-KTGL--TLVL
 bGLUT5 1 MEQDDP-----VKREGRL--TPVI
 hGLUT5 1 MEQDDP-----SMKREGRL--TLVL
 hGLUT7 1 MENKEA-----G--T--P--PPI--PSREGRL--QPTL
 hGLUT1 1 MEPS-------HAT-----FMGKMLLLWSCTSMFPFLGGLSLLVGLV
 hGLUT2 1 MTE-----RGT-----SILPTTLTSLWSLVAIFSVGMIGSFSVGLV
 hGLUT3 1 MGT-----QKV--TPAL
 hGLUT4 1 MPSGFG-----QIGS-----DGE--PPQQRV--TGTL
 HXT7 1 MSQDAI-----AE--QTPVEHLSAVDSASHSVLSTPSNKAERDEIKAYGEGEEHPVVEI--PKRPASA--YVTV
 PfHT1 1 MTKSSKD-----ICS-----E--GKNKSGSGFFTS
 GlcT 1 MOSSTYAVKGNAAFQRRFTSSDRSTSTGIRFAGYKSLATTGPLYCSGSEAMGATLARADNGIQSVMSFSSVKARSVRAQASSDGDDEE--EAIPLRSEGS--SGTV
 XylE 1 MNTQY-----N--SSYI

TM1 56 TM2
 rGLUT5 17 ALATFLMAFGSSFOYGVNVAAVNSPSEFMQQFYNDTYDR-----NKE-----NIESFTTLTLLWSLTVSMFPFGGFIGSLMVGLV
 bGLUT5 18 VLATLTAAPGSSFOYGVNVAAVNSPSEFMKDFNYTYDR-----VGE-----YMEFFYTLTLLWSLTVSMFPFGGFIGSLMVGLV
 hGLUT5 18 ALATLTAAPGSSFOYGVNVAAVNSPALLMQQFYNETYGR-----TGE-----FMEDFPPTLLTLLWSLTVSMFPFGGFIGSLLVGLV
 hGLUT7 24 LLATLTAAPGSAFOYGVNVAAVNSPSEFMKDFNYTYDR-----HAT-----FMGKMLLLWSCTSMFPFLGGLSLLVGLV
 hGLUT1 13 MAVGGAVLGLGNGTGVNAPQKVIEFFYQTFVH-----RGT-----SILPTTLTSLWSLVAIFSVGMIGSFSVGLV
 hGLUT2 11 VFTVTAIVGLG-SFOYGVNVAAVNSPSEFMKDFNYTYDR-----NKE-----NIESFTTLTLLWSLTVSMFPFGGFIGSLMVGLV
 hGLUT3 11 IFAITVATIG-SFOYGVNVAAVNSPSEFMKDFNYTYDR-----VGE-----YMEFFYTLTLLWSLTVSMFPFGGFIGSLMVGLV
 hGLUT4 22 VLAFTVAVGLG-SFOYGVNVAAVNSPSEFMKDFNYTYDR-----TGE-----FMEDFPPTLLTLLWSLTVSMFPFGGFIGSLLVGLV
 HXT7 66 SIMCMTAFG-GFVFGYDGTGTSIGFINQT-DFI-RRFGM-----HAT-----FMGKMLLLWSCTSMFPFLGGLSLLVGLV
 PfHT1 27 FKVYLSMCTA-SFTFGYDGTGTSIGFINQT-DFI-RRFGM-----HAT-----FMGKMLLLWSCTSMFPFLGGLSLLVGLV
 GlcT 106 LPFVGAVLGLG-AIFFGYDGTGTSIGFINQT-DFI-RRFGM-----HAT-----FMGKMLLLWSCTSMFPFLGGLSLLVGLV
 XylE 11 FSITLVATIG-GLIFGDTAVISGTVESLNTVF--VA-----PQ-----NLSAANSLLGFCVASALICIGIIGALGGYCS

TM3 123 TM4 TM5
 rGLUT5 93 NNLEKRGALLFNNIFSILPAILMGCSKI---AK-----SFEIIASRLVGLICAGISSNVVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 bGLUT5 94 NNLEKRGALLFNNIFSILPAILMGCSKI---AK-----SFEIIASRLVGLICAGISSNVVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 hGLUT5 94 NNLEKRGALLFNNIFSILPAILMGCSKI---AK-----SFEIIASRLVGLICAGISSNVVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 hGLUT7 100 DSGRKGALLFNNIFSILPAILMGCSKI---AK-----SFEIIASRLVGLICAGISSNVVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 hGLUT1 88 NRGGRNSMLMNLAFVSAVLMGFSKL---GK-----SFEMLILGRFITIGVCGLTGTFVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 hGLUT2 120 DTGRIKAMLVANLSLVAALLMGFSKL---GP-----SHILIIAGRSISGLVGLICAGISSNVVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 hGLUT3 86 NRGGRNSMLMNLAFVSAVLMGFSKL---GK-----SFEMLILGRFITIGVCGLTGTFVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 hGLUT4 104 QWVGKRNMLVNNVLAFLVGLMGFSKL---GK-----SFEMLILGRFITIGVCGLTGTFVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 HXT7 138 DMVGRKGLIVVVVYIIG-ILIIQIASI---NK-----WQY-FIGRIISGLVGLICAGISSNVVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 PfHT1 100 QFGRRLSLIIYNNFFLVLSILT---SI---TH-----HFHTILFARLSGLVGLICAGISSNVVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 GlcT 169 DKGRGRTFQDLAFLAAGFLCAT---AQ-----SFTQIVMVRGLLAGIGISSNVVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--
 XylE 80 NRGGRNSMLMNLAFVSAVLMGFSKL---GK-----SFEMLILGRFITIGVCGLTGTFVPMYGLGLAPKNLRGALGVVPLFTVGLVAQL--FGL--

TM6 ICH1 ICH2 ICH3 ICH4
 rGLUT5 182 ---RSVLAS-----E-EGPILLGLTGVPAGLQLL--LLPFPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 bGLUT5 183 ---RSLAN-----E-EGPILLGLTGVPAGLQLL--LLPFPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 hGLUT5 183 ---RSLAN-----E-EGPILLGLTGVPAGLQLL--LLPFPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 hGLUT7 189 ---QALGN-----P-AGPILLGLTGVPAGLQLL--LLPFPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 hGLUT1 177 ---DSLMGN-----K-DLPLLLGLTGVPAGLQLL--LLPFPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 hGLUT2 209 ---EFTLGN-----Y-DLPLLLGLTGVPAGLQLL--LLPFPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 hGLUT3 175 ---EFTLGN-----Y-DLPLLLGLTGVPAGLQLL--LLPFPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 hGLUT4 193 ---ESLGT-----A-SLPLLLGLTGVPAGLQLL--LLPFPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 HXT7 227 ---KN-YSN-----S-VGRVPLGLCFAWALFMIG-GMTFPEPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 PfHT1 185 ---AMGEGPKADSTPELTSFALKLRLMFLFPSPVISLGLALVVFPEKEMVFEK-GRIEESKN-ILKKIYETD--NVDPEPLNAKEAVEQNESAKNLSLGLSALK
 GlcT 255 ---P-LAAN-----P-LMRTMFGVAVIPSVLLAI--GMATSPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER
 XylE 188 SGDA-SWLN-----T-DGRYMFASECIPALFLM--LLYTPESPRMMLIQKKNESAAEK-ALQTLRGWK--DVMDEMEERKEDEAEKA--AGFISVVKLFER

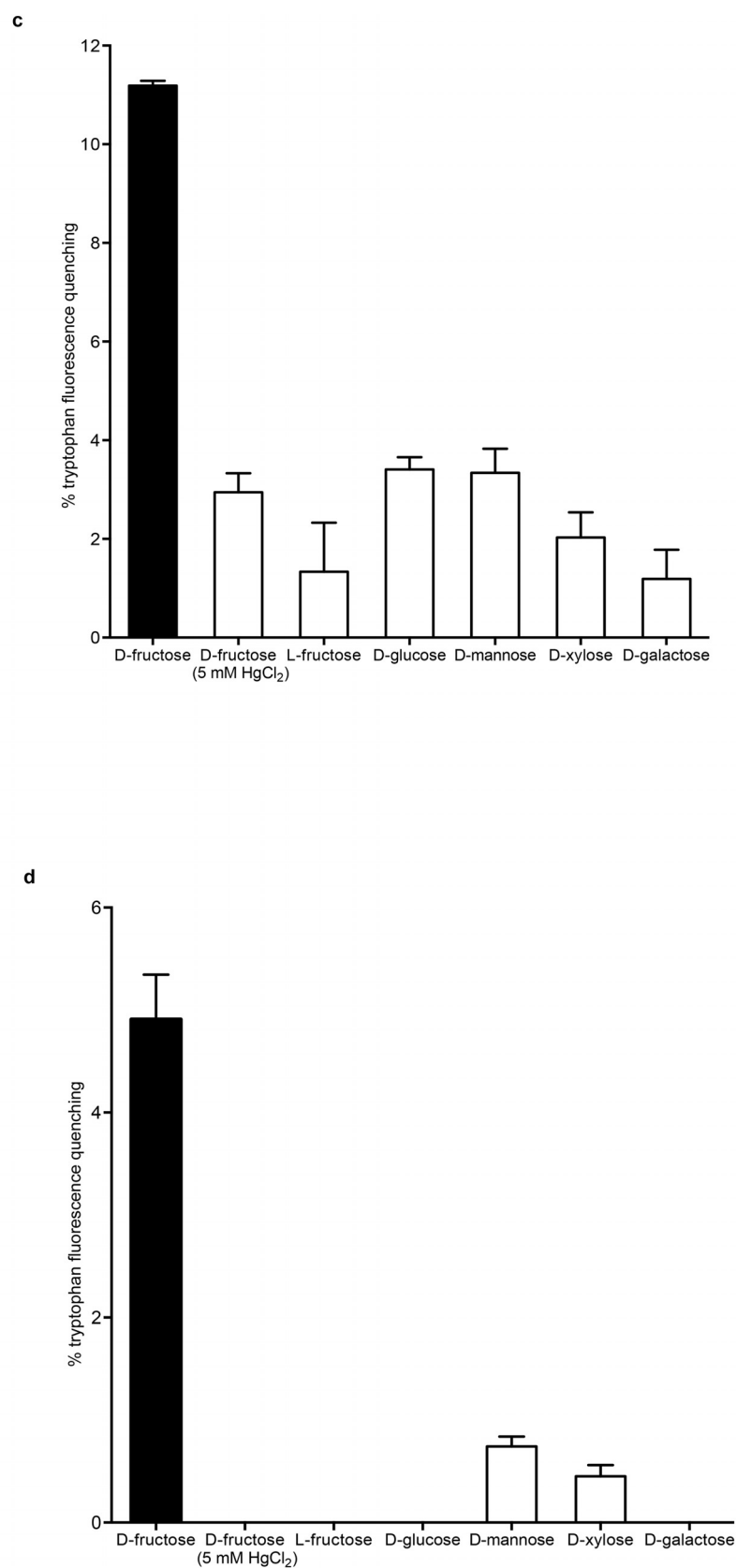
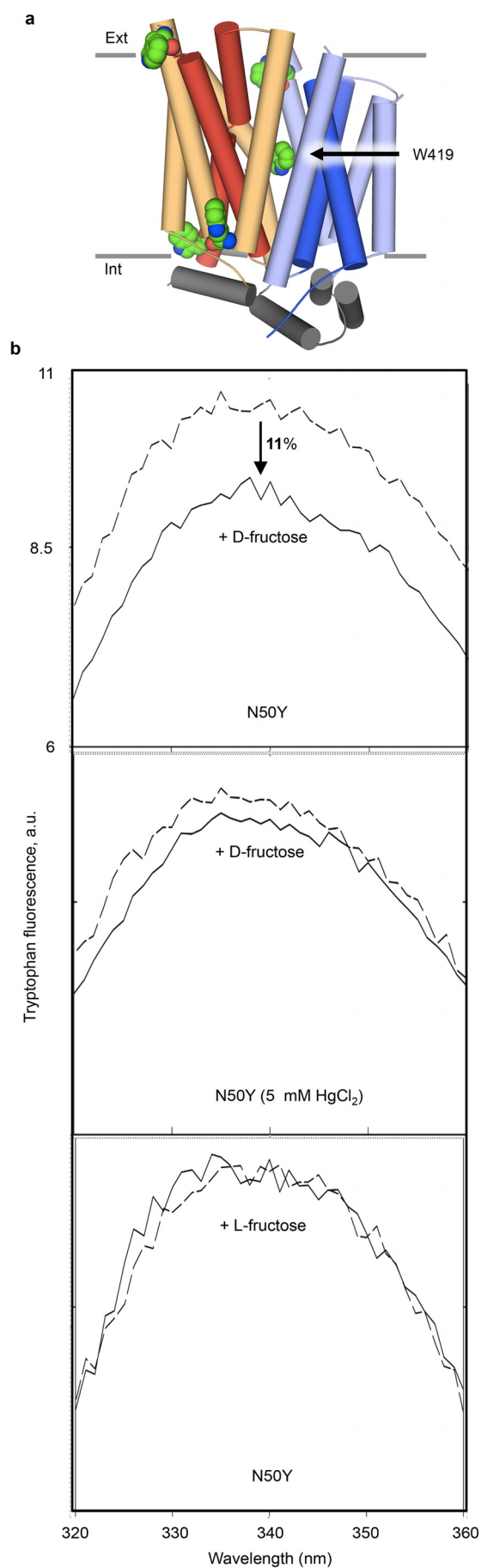
TM7a TM7b TM8 TM9
 rGLUT5 270 MQSL-RWQLISTIVIMTGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 bGLUT5 271 MRSI-RWQVISIIVMAGQGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 hGLUT5 271 MRSI-RWQVISIIVMAGQGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 hGLUT7 277 LRSI-RWQLISTIVIMTGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 hGLUT1 265 SPAY-RQPLIIAIVLQGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 hGLUT2 297 NRSY-RQPLIIAIVLQGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 hGLUT3 263 VSSY-RQPLIIAIVLQGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 hGLUT4 281 SRTH-RQPLIIAIVLQGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 HXT7 317 KSTKVLGRLIMGAMQGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 PfHT1 286 IPSY-RVVIILGCLSGSLGGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 GlcT 340 SR-Y-NKVVSVGAALFLGGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 XylE 275 V-----GVVIGVMTSIFGGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS

TM10 TM11 TM12 ICH5
 rGLUT5 370 WMPYVSVICVIVYVIGVAVGSPDIPALFITIFLQSSRPSNMIGGSGVGLSNFVGLVLPFFIQV-----GLGP-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 bGLUT5 371 WMPYVSVICVIVYVIGVAVGSPDIPALFITIFLQSSRPSNMIGGSGVGLSNFVGLVLPFFIQV-----GLGA-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 hGLUT5 371 WMPYVSVICVIVYVIGVAVGSPDIPALFITIFLQSSRPSNMIGGSGVGLSNFVGLVLPFFIQV-----GLGP-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 hGLUT7 377 ELSYLGICVFAIYAGHSIGPSVPSVVRTIFLQSSRPSNMIGGSGVGLSNFVGLVLPFFIQV-----GLGA-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 hGLUT1 363 WMSYLSIIVAFVGFVAFVFGPDPFVFAVGLSOGQPRPAALAVAGFSNMNTSNFVGLVLPFFIQV-----GLGP-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 hGLUT2 395 WMSYLSIIVAFVGFVAFVFGPDPFVFAVGLSOGQPRPAALAVAGFSNMNTSNFVGLVLPFFIQV-----GLGA-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 hGLUT3 361 GMSVFCIGAILVFAVGFVAFVFGPDPFVFAVGLSOGQPRPAALAVAGFSNMNTSNFVGLVLPFFIQV-----GLGP-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 hGLUT4 379 AMSYLSIIVAFVGFVAFVFGPDPFVFAVGLSOGQPRPAALAVAGFSNMNTSNFVGLVLPFFIQV-----GLGA-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 HXT7 424 GAGNCMIVFACFYIFCPTATWADIPYVVSSTPLRVKSKMSASLATAANLWGLFGLFFPFIQV-----GLGP-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 PfHT1 387 FVKILSIIVAFVGFVAFVFGPDPFVFAVGLSOGQPRPAALAVAGFSNMNTSNFVGLVLPFFIQV-----GLGA-YSFIFVAVICLLTSIYFMVVPVKKRTFVEIN
 GlcT 437 YSGTLAVVGTVLYVLSGLGGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS
 XylE 369 ---GIVALLSMLFYVAAVAMSGLVGLSGLVNAIYVADQIYLSAAGVKSNDVQVYTAGTAVVFMVMTVFVVLWGRNMLLIQFSTCLTACIVLTVALALQN-----TIS

rGLUT5 472 QIFAKNMKV-----SDV-----YPEKEEKELN--DLP-----PATREQ
 bGLUT5 473 RIFIKNMKV-----PGV-----HPEKE--ELK--EFP-----PSTARQ
 hGLUT5 473 QIFAKNMKV-----SEV-----YPEKE--ELK--ELP-----PVTSQ
 hGLUT7 479 RIFAKNMKV-----KL-----PEEKEE-TID--AGPPTASPAKETS
 hGLUT1 465 SGFRQGA-----SQS-----DKTPE--ELF--PLGADSOV
 hGLUT2 497 AEFQKSGS-----AHR-----PKAA--VEMK-----FLGATETV
 hGLUT3 463 RAFEGQAHG-----ADR-----SGKDGVMEMN--SIEP--AKETTNTV
 hGLUT4 481 AAFHRTPSL-----LEQ-----EVKPS--TELE--YLPDEND
 HXT7 526 TMEEGVLWPKSASVPPSRRGANYDAE--EMTHDDKPL--YKRMFSK
 PfHT1 489 YITM-----YLM-----ERQ-----KMTKSVV
 GlcT 540 LAL-----TSQA
 XylE 476 ALWEPETKK-----TQQTATL

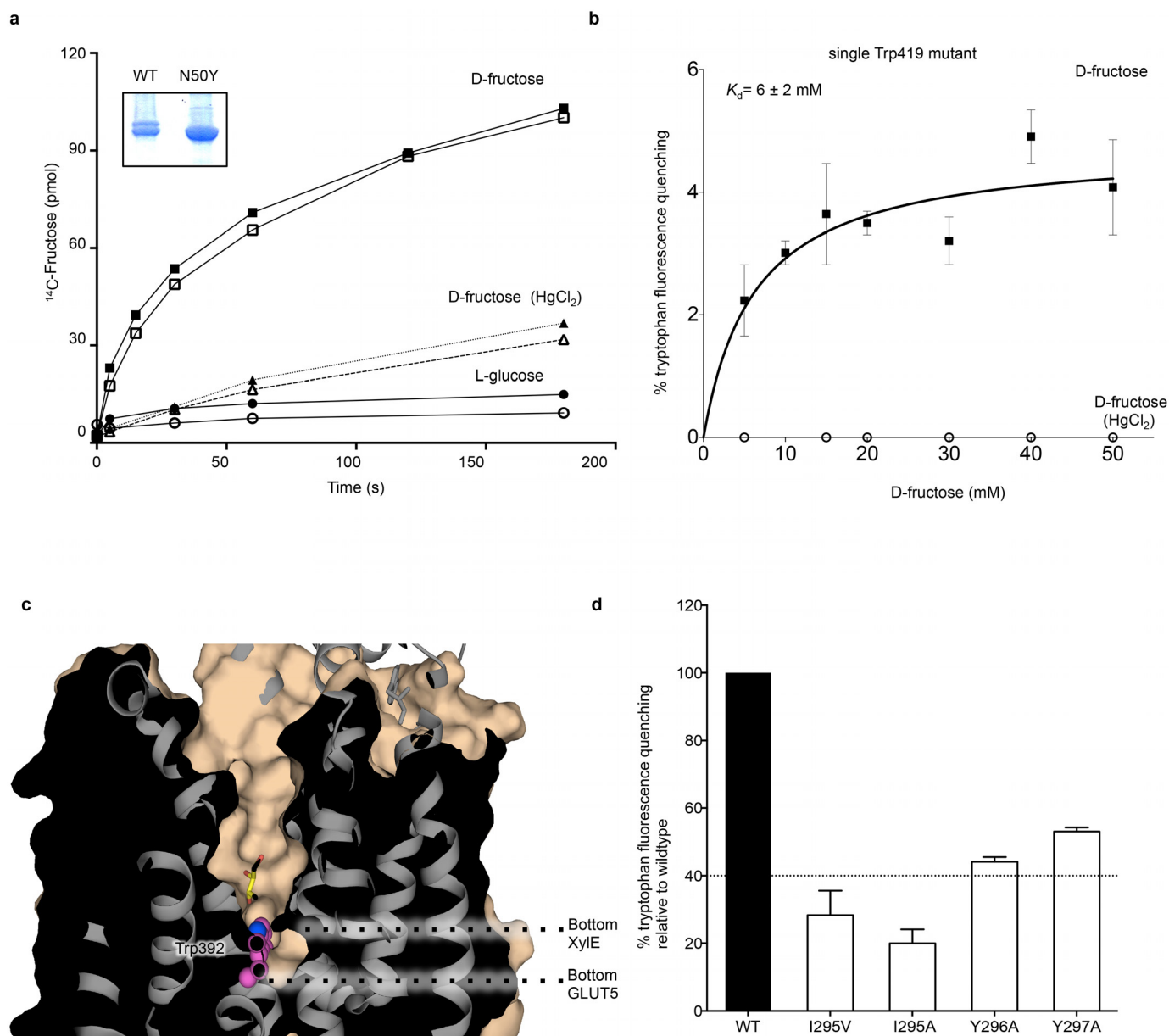
Extended Data Figure 4 | Sequence alignment of rat GLUT5 (rGLUT5), bovine GLUT5 (bGLUT5), human GLUT5 and GLUT7 (hGLUT5 and hGLUT7), human GLUT1–4 (hGLUT1–4), *Saccharomyces cerevisiae* HXT7, *Plasmodium falciparum* PfHT, *Arabidopsis thaliana* GlcT and *Escherichia coli* XylE. Structure elements of rat GLUT5 are indicated above the alignment, and coloured as in Fig. 1a. Strictly conserved residues are highlighted in black-filled boxes, and highly conserved residues are shaded in grey. Green boxes highlight central cavity residues that are specific to GLUT5 and red boxes highlight those that are conserved among GLUTs. Purple boxes highlight residues forming the salt bridges between cytosolic TM segments. A

blue box (TM5) highlights Gln166, whose mutation to glutamic acid, as present in GLUT7, weakens D-fructose binding but supports strong D-glucose binding in rGLUT5. The brown box (TM8) highlights Glu336 that is conserved across all the GLUTs and replaced with glutamic acid in XylE. Red bars underneath the alignment indicate the sugar porter (SP) family motifs^{18,19}. Note that because bGLUT5 and hGLUT5 have an additional amino acid at position 8, their numbering differs from rGLUT5 by 1 amino acid. For clarity, bGLUT5 residues are labelled using rGLUT5 numbering.



Extended Data Figure 5 | D-fructose binding monitored by tryptophan fluorescence quenching. **a**, Cartoon representation of the outward-facing rGLUT5 structure, as viewed from the plane of the membrane with the colouring as shown in Fig. 1a. Atoms in all tryptophan residues are shown as spheres and tryptophan W419, whose fluorescence is quenched by substrate, is labelled. **b**, Emission fluorescence spectra for purified deglycosylated rGLUT5 wild-type-like mutant N50Y (referred to as WT), shown in the range of 320–360 nm with an excitation wavelength of 295 nm after the addition of 40 mM D-fructose (top), and 40 mM L-fructose (bottom). Emission fluorescence spectra for purified wild-type protein that had been previously incubated with the inhibitor HgCl₂ is also shown for D-fructose (middle). **c**, Tryptophan

fluorescence quenching (excitation 295 nm; emission 338 nm) after incubation of purified rGLUT5 N50Y with either 40 mM D-fructose (filled bar) or L-fructose, D-glucose, D-mannose, D-xylose or D-galactose as labelled (open bars). Tryptophan fluorescence quenching for purified wild-type protein that had been previously incubated with the inhibitor HgCl₂ is also shown for D-fructose (open bar). **d**, As in **c**, but for rGLUT5 with a single tryptophan residue (W419), which contains the following mutations: N50Y, W70F, W191F, W239F, W265F, W275F, W338F and W370F. No tryptophan quenching was observed for D-fructose (5 mM HgCl₂), L-fructose, D-glucose or D-galactose. In all experiments errors bars indicate s.e.m.; *n* = 3.

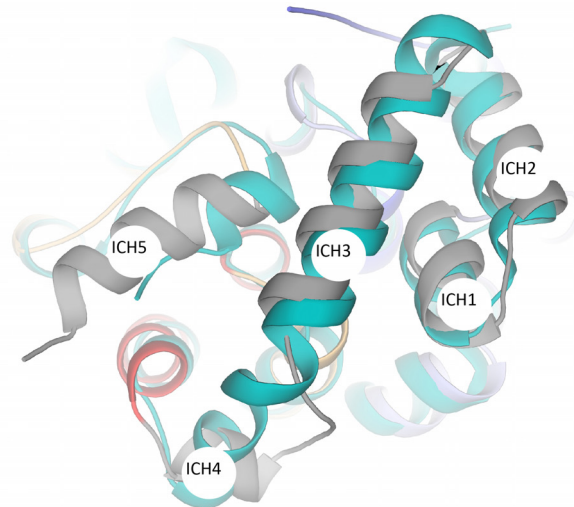
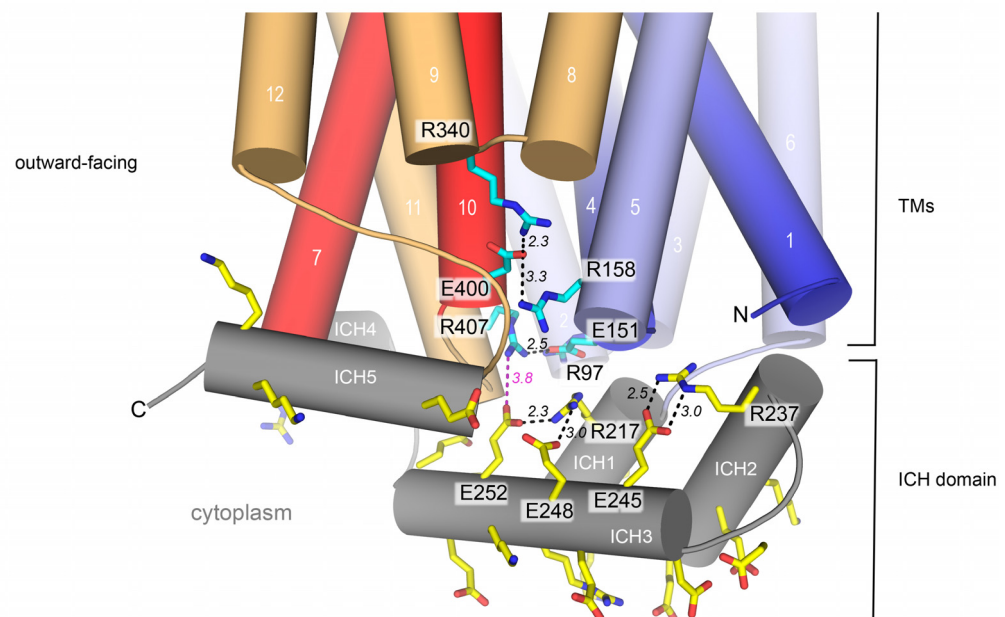
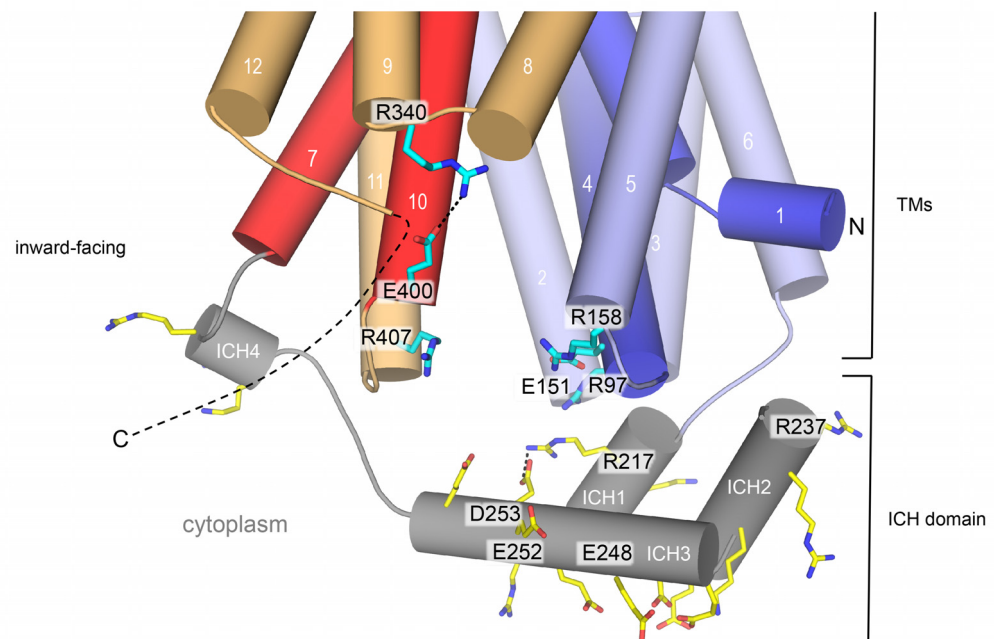


Extended Data Figure 6 | Substrate specificity in GLUT5. **a**, Time-dependent uptake of D- ^{14}C -fructose by rGLUT5 wild type (open squares and triangles) and the deglycosylated mutant N50Y (filled squares and triangles) in proteoliposomes incubated with or without the inhibitor HgCl_2 as labelled. Non-specific uptake was estimated with 0.1 mM L- ^{14}C -glucose for wild type (filled circles) and the N50Y mutant (open circles). In all experiments errors bars represent a spread of duplicates. Inset shows SDS-PAGE analysis of the purified rat GLUT5 wild type and the deglycosylated N50Y mutant. **b**, Tryptophan fluorescence quenching (excitation 295 nm; emission 338 nm), after incubation of purified rat GLUT5 mutant (N50Y, W70F, W191F, W239F, W265F, W275F, W338F, W370F) that contains one single tryptophan residue, W419, with increasing concentrations of D-fructose (filled squares) and to the protein previously incubated with the inhibitor mercury chloride

(open circles). **c**, Slab through the surface of the outward-facing rGLUT5 structure as viewed in the plane of membrane. The structure of substrate-bound XylE structure was further superimposed onto rGLUT5 and is shown here as a grey ribbon. In XylE, Trp392 (Trp388 in hGLUT1) is located at the bottom of the cavity (spheres; magenta) and coordinates D-xylose (stick form; yellow). In GLUT5, the equivalent residue is an alanine, making the cavity deeper. **d**, D-fructose binding as measured by tryptophan fluorescence quenching (excitation 295 nm; emission 338 nm) after incubation with 40 mM D-fructose for wild type (open bar), and TM7 mutations of Ile295 (interacts with TM10 residues) and Tyr296 and Tyr297 residues. Equivalently located tyrosine residues in XylE occlude the sugar-binding site from the outside²². Fluorescence quenching for the mutants are displayed as a percentage of total wild-type binding. In all experiments errors bars indicate s.e.m.; $n = 3$.

a

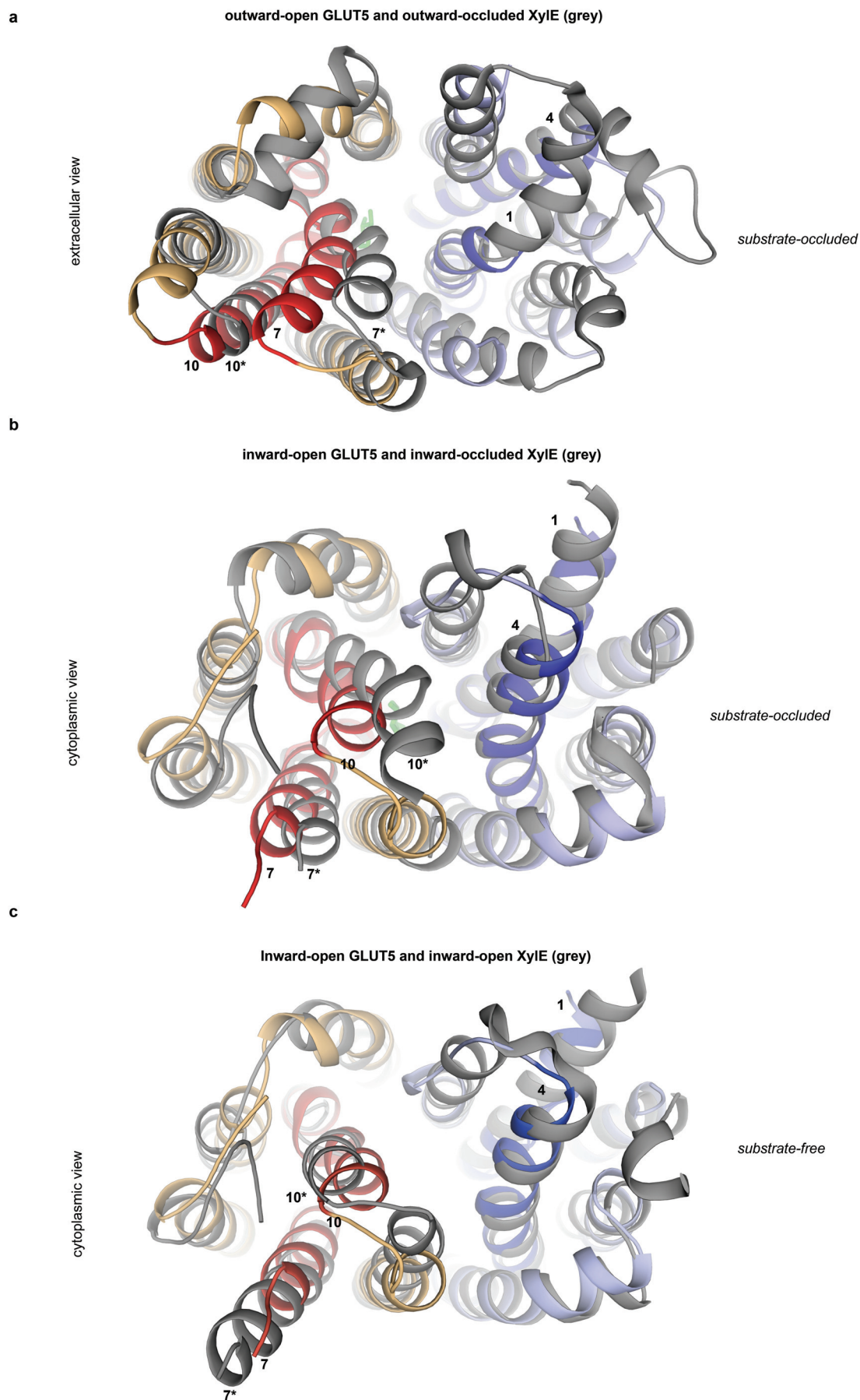
grey = outward-facing Glut5
teal = outward-occluded XylE

**b****c**

Extended Data Figure 7 | The intracellular helical domain (ICH).

a, Cytoplasmic view of the ICH domain after superposition of the open, outward-facing rGLUT5-sFv (grey) and outward-facing occluded *E. coli* Xyle (teal) (4GBY) structures. **b**, In the outward-facing GLUT5 structure ICH1–ICH3 are linked together by several salt bridges (side chains are labelled and shown as sticks in yellow). In contrast, no polar interactions are formed between ICH5 and either ICH1–ICH3 or cytoplasmic ends of N-terminal TM bundle helices. A salt bridge forms (dotted line in magenta), however, between

Glu225 in ICH3 and Arg407 in TM11, which also forms part of the inter-bundle salt-bridge network (side chains are labelled and shown as sticks in cyan). **c**, In the inward-facing GLUT5 structure, this inter-bundle salt-bridge network is not formed, because the cytoplasmic ends of the N- and C-terminal bundle have moved apart; consistently, the ICH domain functional role is proposed to act as a scaffold domain that further helps to stabilize the outward-facing conformation²¹.



Extended Data Figure 8 | Access to the central cavity and substrate-binding site is gated by TM7 on the outside and TM10 on the inside. **a**, Superposition of outward-facing open GLUT5 and outward-facing occluded *E. coli* XylE (4GBY) structures. The TM numbering for outward-facing occluded XylE has an additional asterisk. The inward-facing GLUT5 structure is coloured as in Fig. 1a and that of XylE in grey. The bound D-xylose is shown in stick representation in green. The r.m.s.d. is 1.38 Å for 290 pairs of C α atoms (see Methods). **b**, Superposition of inward-open GLUT5 and inward-occluded

E. coli XylE structure (4JA3) with colouring and annotation as described in **a**. The r.m.s.d. is 1.80 Å for 274 pairs of C α atoms (see Methods). The bound D-xylose in 4GBY is represented in stick form in green. The ICH domain is not shown for clarity. **c**, Superposition of inward-facing open GLUT5 and inward-facing open XylE (4JA4) structures as viewed from the cytoplasmic side with colouring and annotation as described in **a**. The ICH domain is not shown for clarity. The r.m.s.d. is 1.70 Å for 273 pairs of C α atoms (see Methods).

Extended Data Table 1 | Crystallographic data collection and refinement statistics

	rat GLUT5-Fv	bovine GLUT5 [#]
Data collection		
Space group	$P2_1$	$P2_12_12_1$
Cell dimensions		
<i>a</i> , <i>b</i> , <i>c</i> (Å)	76.78, 151.54, 106.40	74.61, 112.15, 139.57
α , β , γ (°)	90.00, 97.25, 90.00	90.00, 90.00, 90.00
Resolution (Å)	50-3.27 (3.39-3.27)	100 - 3.00 (3.11 - 3.01)*
R_{sym} or R_{merge}	21.5 (>100)	10.06 (>100)
<i>I</i> / σ <i>I</i>	11.6 (1.5)	11.21 (0.86)
Completeness (%)	100 (100.0)	99.4 (93.9)
Redundancy	13.0 (13.0)	10.2 (8.6)
Refinement		
Resolution (Å)	50-3.27 (3.5-3.27)	33.3 - 3.2 (3.31 - 3.20)
No. reflections (Rfree set)	38017 (3355)	13346 (1331)
R_{work} / R_{free}	24.2/28.8 (35.4/37.5)	23.6/25.8 (32.8/36.5)
No. atoms		
Protein	10657	3382
B-factors		
Protein	157.0	149.2
R.m.s deviations		
Bond lengths (Å)	0.004	0.003
Bond angles (°)	0.97	0.90
Ramachandran plot (%)		
Favored	98.0	94.5
Outliers	0	0

[#] Data were obtained by scaling together two datasets collected on the same crystal

*Highest resolution shell used in the final refinement is shown in parenthesis.

Extended Data Table 2 | Completeness of bovine GLUT5 data per resolution shell after correction for anisotropy

Resolution range (Å)	Completeness (%)
100.0 - 8.10	95.5
8.10 - 6.43	100.0
6.43 - 5.62	99.9
5.62 - 5.10	99.9
5.10 - 4.74	100.0
4.74 - 4.46	99.7
4.46 - 4.24	99.6
4.24 - 4.05	100.0
4.05 - 3.90	78.4
3.90 - 3.76	53.4
3.76 - 3.64	40.6
3.64 - 3.54	33.4
3.54 - 3.45	26.3
3.45 - 3.36	21.2
3.36 - 3.29	15.4
3.29 - 3.22	13.1
3.22 - 3.15	9.4
3.15 - 3.09	6.0
3.09 - 3.04	3.5
3.04 - 3.00	1.7

Two independent and primitive envelopes of the bilobate nucleus of comet 67P

Matteo Massironi^{1,2}, Emanuele Simioni³, Francesco Marzari⁴, Gabriele Cremonese⁵, Lorenza Giacomini¹, Maurizio Pajola², Laurent Jorda⁶, Giampiero Naletto^{2,3,7}, Stephen Lowry⁸, Mohamed Ramy El-Maarry⁹, Frank Preusker¹⁰, Frank Scholten¹⁰, Holger Sierks¹¹, Cesare Barbieri⁴, Philippe Lamy⁶, Rafael Rodrigo^{12,13}, Detlef Koschny¹⁴, Hans Rickman^{15,16}, Horst Uwe Keller¹⁷, Michael F. A'Hearn^{18,19}, Jessica Agarwal¹¹, Anne-Thérèse Auger⁶, M. Antonella Barucci²⁰, Jean-Loup Bertaux²¹, Ivano Bertini², Sebastien Besse¹⁴, Dennis Bodewits¹⁸, Claire Capanna⁶, Vania Da Deppo³, Björn Davidsson²², Stefano Debei²³, Mariolino De Cecco²⁴, Francesca Ferri², Sonia Fornasier²⁰, Marco Fulle²⁵, Robert Gaskell²⁶, Olivier Groussin⁶, Pedro J. Gutiérrez²⁷, Carsten Güttler¹¹, Stubbe F. Hviid^{10,11}, Wing-Huen Ip²⁸, Jörg Knollenberg¹⁰, Gabor Kovacs¹¹, Rainer Kramm¹¹, Ekkehard Kühr¹⁰, Michael Küppers²⁹, Fiorangela La Forgia⁴, Luisa M. Lara²⁷, Monica Lazzarin⁴, Zhong-Yi Lin²⁸, José J. Lopez Moreno²⁷, Sara Magrin⁴, Harald Michalik³⁰, Stefano Mottola¹⁰, Nilda Oklay¹¹, Antoine Pommerol⁹, Nicolas Thomas⁹, Cecilia Tubiana¹¹ & Jean-Baptiste Vincent¹¹

The factors shaping cometary nuclei are still largely unknown, but could be the result of concurrent effects of evolutionary^{1,2} and primordial processes^{3,4}. The peculiar bilobed shape of comet 67P/Churyumov–Gerasimenko may be the result of the fusion of two objects that were once separate or the result of a localized excavation by outgassing at the interface between the two lobes⁵. Here we report that the comet's major lobe is enveloped by a nearly continuous set of strata, up to 650 metres thick, which are independent of an analogous stratified envelope on the minor lobe. Gravity vectors computed for the two lobes separately are closer to perpendicular to the strata than those calculated for the entire nucleus and adjacent to the neck separating the two lobes. Therefore comet 67P/Churyumov–Gerasimenko is an accreted body of two distinct objects with 'onion-like' stratification, which formed before they merged. We conclude that gentle, low-velocity collisions occurred between two fully formed kilometre-sized cometesimals in the early stages of the Solar System. The notable structural similarities between the two lobes of comet 67P/Churyumov–Gerasimenko indicate that the early-forming cometesimals experienced similar primordial stratified accretion, even though they formed independently.

Soon after the insertion of the robotic space probe Rosetta into orbit around comet 67P/Churyumov–Gerasimenko on 6 August 2014, images acquired by the Optical, Spectroscopic, and Infrared Remote Imaging System (OSIRIS)⁶ have provided evidence of morphological terraces distributed all over its surface^{5,7}. In the following months, imaging data from both the Narrow Angle Camera (NAC) and Wide Angle Camera (WAC) have confirmed that these features, together with cuesta-like morphologies, are related to a pervasive stratification at the surface of the nucleus (Figs 1, 2 and Extended Data Figs 1–4). (A cuesta is an asymmetric ridge characterized by a long and gentle backslope that conforms with the dip of a resistant stratum called cap-rock, and a steep

frontal slope.) Terraces and localized strata are not new features of cometary nuclei, having been observed at the surface of comet 9P/Tempel 1 (refs 8 and 9) and, possibly, of comets 81P/Wild 2 (ref. 10) and 19P/Borrelly (ref. 11). Their origin has been explained either as the result of sublimation, erosion and redeposition of lag materials at the surface¹, as successive flows of fluidized material², or as the surface expression of an inner stratified structure related to primordial accretion^{3,4}. However, the OSIRIS images of comet 67P/Churyumov–Gerasimenko provide evidence of global stratification with an imaging resolution of up to 0.1 m per pixel, greatly exceeding that of previous close-range images, the best of which are around 7 m per pixel^{12,13} (Figs 1, 2 and Methods section 'Data and surface strata description').

The structure delineated by the strata and the relationship between gravity vectors and strata orientations provide information that we can use to disentangle the main formation hypotheses concerning the unusual bilobed shape of the nucleus of comet 67P/Churyumov–Gerasimenko. Is its shape the manifestation of a concentrated localized erosion of a single monolithic body into a two-lobe configuration or the expression of two distinct accreted objects? What are the implications in terms of comet origin?

We were able to retrieve the strata orientations of both the major lobe (the main body) and the minor lobe (the head) of the comet by using the best-fitting planes of terraces and cuesta-like features obtained from a stereo-photoclinometry shape model covering 75% of the comet nucleus^{5,7} (Methods section 'Best-fitting planes'). The 103 best-fitting planes throughout the entire comet nucleus that we derived emphasize that the strata are coherently oriented all around the comet, suggesting two separate wrapping sequences for the head and the main body (Fig. 3 and Extended Data Figs 5 and 6).

The strata are ubiquitous at the surface, vertically outspread along scarps, parallel to each other, and laterally continuous (Methods

¹Dipartimento di Geoscienze, University of Padova, via G. Gradenigo 6, 35131 Padova, Italy. ²Centro di Ateneo di Studi ed Attività Spaziali "Giuseppe Colombo" (CISAS), University of Padova, via Venezia 15, 35131 Padova, Italy. ³CNR-IFN UOS Padova LUXOR, via Trasea 7, 35131 Padova, Italy. ⁴University of Padova, Department of Physics and Astronomy, Vicolo dell'Osservatorio 3, 35122 Padova, Italy.

⁵INAF, Osservatorio Astronomico di Padova, Vicolo dell'Osservatorio 5, 35122 Padova, Italy. ⁶Aix Marseille Université, CNRS, LAM (Laboratoire d'Astrophysique de Marseille), UMR 7326, 38 rue Frédéric Joliot-Curie, 13388 Marseille, France. ⁷Department of Information Engineering, University of Padova, via Gradenigo 6/B, 35131 Padova, Italy. ⁸The University of Kent, School of Physical Sciences, Canterbury, Kent CT2 7NZ, UK. ⁹Physikalisches Institut der Universität Bern, Sidlerstraße 5, 3012 Bern, Switzerland. ¹⁰Deutsches Zentrum für Luft- und Raumfahrt (DLR), Institut für Planetenforschung, Rutherfordstraße 2, 12489 Berlin, Germany. ¹¹Max-Planck-Institut für Sonnensystemforschung, Justus-von-Liebig Weg 3, 37077 Göttingen, Germany. ¹²Centro de Astrobiología, CSIC-INTA, 28850 Torrejón de Ardoz, Madrid, Spain. ¹³International Space Science Institute, Hallerstraße 6, 3012 Bern, Switzerland. ¹⁴Scientific Support Office, European Space Research and Technology Centre/ESA, Keplerlaan 1, Postbus 299, 2201 AZ Noordwijk ZH, The Netherlands. ¹⁵Department of Physics and Astronomy, Uppsala University, Box 516, 75120 Uppsala, Sweden. ¹⁶PAS Space Research Center, Bartycka 18A, 00716 Warszawa, Poland. ¹⁷Institut für Geophysik und extraterrestrische Physik (IGEP), Technische Universität Braunschweig, Mendelssohnstraße 3, 38106 Braunschweig, Germany.

¹⁸University of Maryland, Department of Astronomy, College Park, Maryland 20742-2421, USA. ¹⁹Akademie der Wissenschaften zu Göttingen and Max-Planck-Institut für Sonnensystemforschung, Justus-von-Liebig Weg 3, 37077 Göttingen, Germany. ²⁰LESIA-Observatoire de Paris, CNRS, Université Pierre et Marie Curie, Université Paris Diderot, 5 place J. Janssen, 92195 Meudon, France. ²¹LATMOS, CNRS/UVSQ/IPSL, 11 boulevard d'Alembert, 78280 Guyancourt, France. ²²Department of Physics and Astronomy, Uppsala University, Box 516, 75120 Uppsala, Sweden. ²³Department of Industrial Engineering, University of Padova, via Venezia 1, 35131 Padova, Italy. ²⁴University of Trento, Via Mesiano 77, 38100 Trento, Italy. ²⁵INAF—Osservatorio Astronomico, Via Tiepolo 11, 34014 Trieste, Italy.

²⁶Planetary Science Institute, Tucson, Arizona 85719, USA. ²⁷Instituto de Astrofísica de Andalucía (CSIC), Glorieta de la Astronomía s/n, 18008 Granada, Spain. ²⁸National Central University, Graduate Institute of Astronomy, 300 Chung-Da Road, Chung-Li 32054 Taiwan. ²⁹Operations Department, European Space Astronomy Centre/ESA, PO Box 78, 28691 Villanueva de la Canada, Madrid, Spain.

³⁰Institut für Datentechnik und Kommunikationsnetze der TU Braunschweig, Hans-Sommer Straße 66, 38106 Braunschweig, Germany.

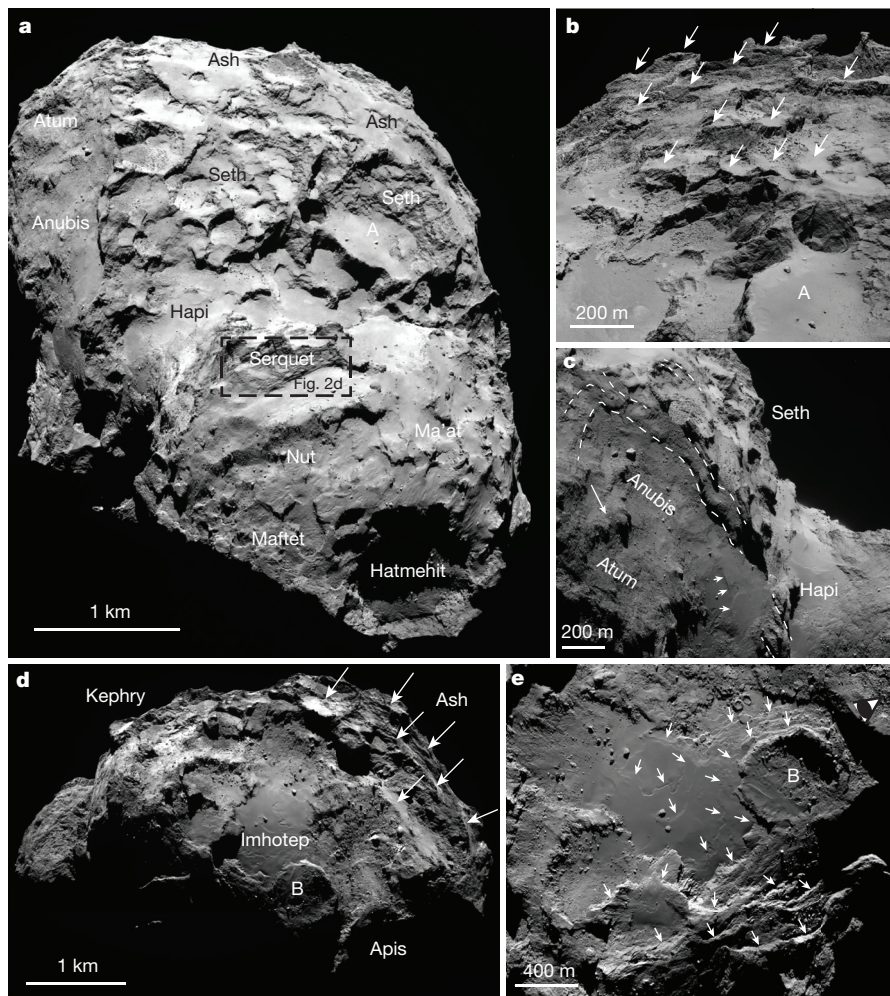


Figure 1 | The main body. **a**, Overview of the Seth region, showing lineated cliffs and terraces. Other regions and the location of the candidate landing site A (one of the five locations identified in August 2014 as candidates to set down the Philae lander) are also shown. The dashed square is the location of Fig. 2d. **b**, Lateral view of site A with Seth terraces in the background (white arrows). Notice the strata within the pit at the margin of site A. **c**, Seth–Anubis contact. Dashed white lines are strata. The long white arrow is a terrace in the Atum region. Small white arrows indicate the terrace margin in Anubis. **d**, WAC view of the Imhotep region. The plateau structure B and the stratification of bounding regions (white arrows) are evident. **e**, NAC view of Imhotep staircase terraces (white arrows) extending underneath the Imhotep southern boundary strata. The eye symbol indicates the view of Extended Data Fig. 2d.

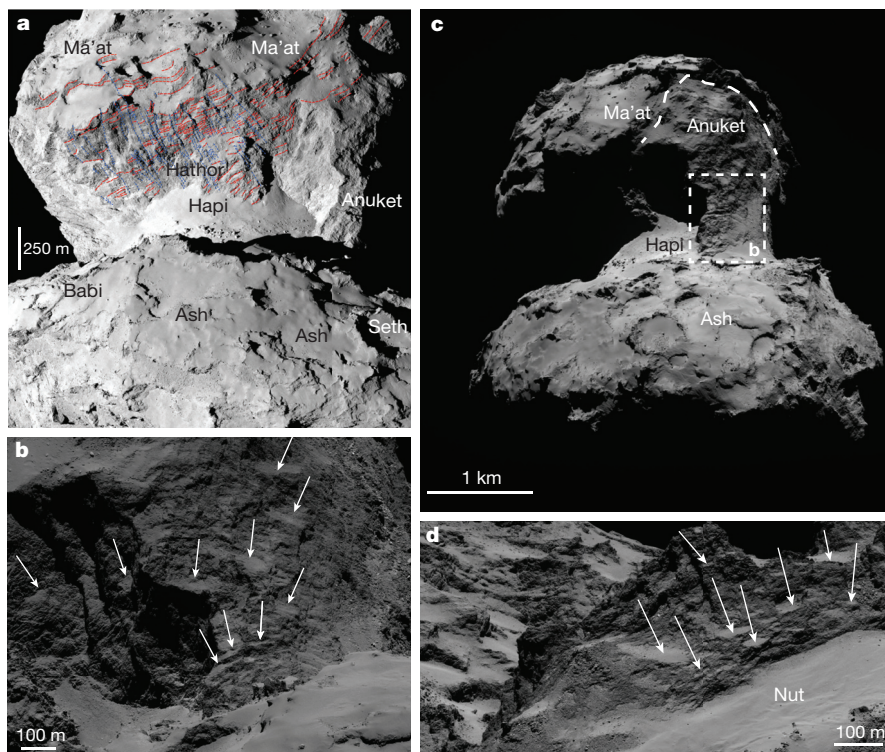


Figure 2 | The head of comet 67P/Churyumov-Gerasimenko. **a**, The Hathor region, showing strata (dashed red lines) and fractures (dashed blue lines). In the foreground the cuesta-like morphologies of the Ash and Babi regions are visible. No abrupt change of strata orientation occurs between Hathor and Ma'at. **b**, Stratification at the Hathor–Anuket boundary. Terraces (white arrows) show that strata are continuously changing their orientation. **c**, WAC overview of Anuket. The dashed white line underlines the scarp where Anuket extends underneath Ma'at. The dashed box is the location of **d**. **d**, NAC view of the Serquet and Nut regions: Serquet terraces (white arrows) are parallel to the Nut planar region.

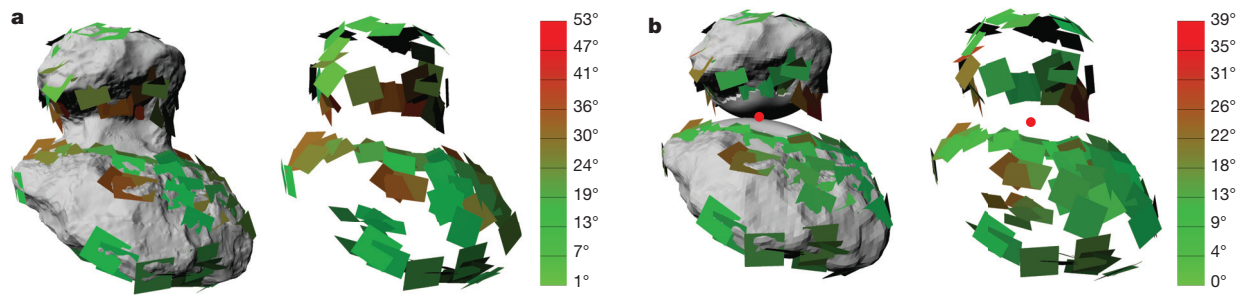


Figure 3 | 3D views of the nucleus and strata of comet 67P/Churyumov-Gerasimenko. **a**, Entire nucleus; **b**, separated lobes. Best-fitting planes derived from strata are shown for the complete shape (left) and alone (right). The

colour scales indicate the angular deviation between the plane vector and the local gravity vector calculated for the whole body (**a**) and the separated lobes (**b**). The red point in **b** is the contact point between the two lobes.

section ‘Surface strata description’), so we were able to obtain a series of geological sections by cutting the comet nucleus along two different sets of planes: one along the major axis and one perpendicular to it (Methods subsection ‘Geological sections’). The geological sections transecting both the main body and the head show that the strata of the two lobes are clearly independent of each other, and define two envelopes not only on the surface but underneath, locally reaching a depth of 650 m (Fig. 4a and Extended Data Figs 7 and 8). The separation is particularly evident in the opposite inclination of strata at the neck flanks (Fig. 4a and Extended Data Figs 7 and 8). Similarly, strata clearly define a mantling envelope on the sections perpendicular to the major axis of the body (Fig. 4b, c and Extended Data Fig. 9).

The thicknesses of the stratified sequence (several hundred metres beneath most regions) we derived should be considered lower bounds because the stratification might continue to greater depths. The thickness of the stratified rocky-like material excludes by itself any causative process related to weathering during the evolution of comet 67P/Churyumov-Gerasimenko. Layering from metamorphic waves should also be ruled out, since metamorphic fronts generate gradual

transitions and not the distinct strata with sharp boundaries that are seen on comet 67P/Churyumov-Gerasimenko.

We can constrain the structural reconstruction further by comparing the orientation of the planes fitting terraces and cuesta-like features with the local gravity-field vectors computed for either the whole comet nucleus or for the two lobes as if they are independent (Methods section ‘Gravity-field vectors’). In Figs 3 and 5 the angular relationship between the planes of the strata and the gravity-field vectors are visualized through three-dimensional (3D) reconstructions and diagrams. The histograms of Fig. 5 show the distribution of the absolute angular deviations between the local gravity vector and the planes normal to the strata at distances lower and higher than 1.7 km from the contact point of the two separated (and reconstructed) lobes (red point in Fig. 3b). The angular values and the standard deviations of the strata close to the neck (distances ≤ 1.7 km) are higher if the gravity field of the entire comet nucleus is considered with respect to the two separated objects. In contrast, no major differences are recorded between the two gravity models far from the neck region (distances >1.7 km), where both angular values and standard deviations are low.

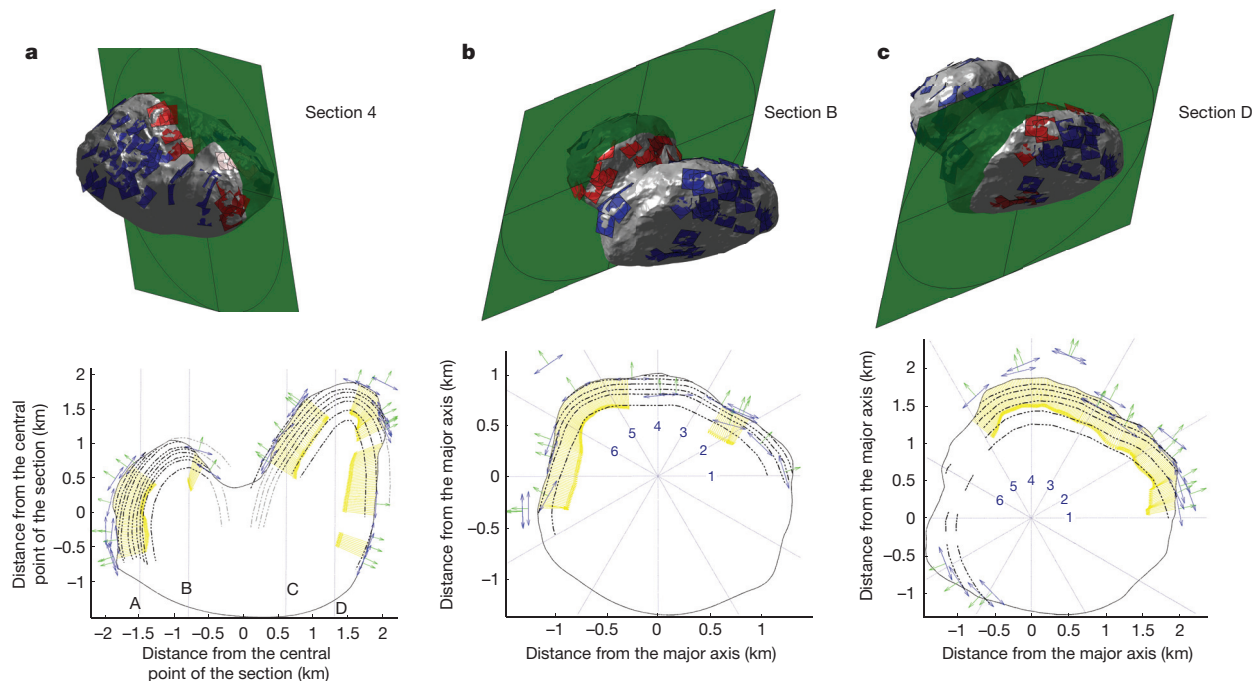


Figure 4 | Geological sections through comet 67P/Churyumov-Gerasimenko. **a**, Longitudinal section along the comet’s major axis ($y = 0$). **b**, **c**, Transversal sections perpendicular to the comet’s major axis ($x = 0$, $y = 0$). In the top row are perspective views of the comet nucleus showing the different cross-sectional slices. Shown in red are the best-fitting planes projected into the geological section; shown in blue are all the other best-fitting planes. In the

bottom row are geological sections. Shown in blue are the best-fitting planes; green arrows are vectors normal to each plane; yellow shading is the field of lines used for drawing strata; dashed black lines indicate strata; dashed grey lines indicate inferred strata. Further geological sections are in Extended Data Figs 7–9. The traces of transversal sections are reported in the longitudinal one and vice versa.

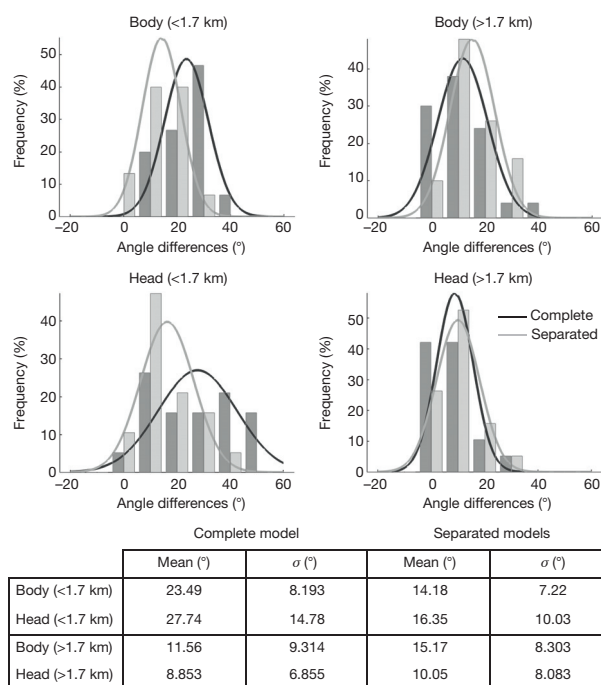


Figure 5 | Strata and local gravity vectors. Distribution of angular deviations between the local gravity vector and the normal-to-strata planes at distances lower and higher than 1.7 km from the point of intersection between the two reconstructed lobes (red point in Fig. 3b). The diagrams show histograms of the angular deviations as a percentage of the total best-fitting planes at distances lower and higher than 1.7 km for the body and the head respectively. Gaussian distributions normalized to the bin dimensions are overplotted on the histograms to emphasize the different means and standard deviations σ (also reported in the table). The angular values are referred to the local gravity vector calculated for the complete comet nucleus (dark grey bars and black curve) as well as for the two separated objects (pale grey bars and grey curve).

Hence the detected stratification seems to be best modelled as being perpendicular to the gravity vectors in each of the lobes. This would suggest that the stratified envelopes of the head and the main body formed independently and before the aggregation of the two lobes.

OSIRIS data have revealed that stratification is a dominant structural aspect of comet 67P/Churyumov-Gerasimenko throughout its entire bilobed shape, with the only exception being the neck region separating the main body from the head. The structural reconstruction of the comet nucleus based on the orientation and distribution of the identified strata shows that the two lobes were two independent objects—perhaps formed as pebble-pile planetesimals consisting of aggregates of primordial pebbles from the Solar nebula¹⁴—characterized by onion-like stratification several hundred metres thick. The cuesta-like morphologies, derived from differential sublimation of the stratified sequences (Methods section ‘Surface strata description’), imply variations in the relative abundances of volatile materials among strata. The relationship between strata plane orientations and gravity-field vectors, calculated for the entire comet and the two separated (and reconstructed) lobes, suggests that these envelopes formed independently before the fusion of the two objects (the head and the main body) into a single consolidated nucleus. The ordered sequence of wrapping strata precludes a chaotic structure up to a depth of several hundred metres, suggests that stratification is a primary structure, and implies a merging of two fully formed kilometre-sized cometesimals in the early stages of the Solar System via a low-velocity impact¹⁵. The two lobes are strikingly similar in terms of deep onion-like structure, surface composition¹⁶ and observed surface features⁷. Taken together, the structural and compositional similarities indicate that the two cometesimals experienced similar accretion processes, having formed completely independently.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 2 March 2015; accepted 10 July 2015.

Published online 28 September 2015.

- Basilevsky, T. & Keller, H. U. Comet nuclei: morphology and implied processes of surface modification. *Planet. Space Sci.* **54**, 808–829 (2006).
- Belton, M. J. S. & Melosh, J. Fluidization and multiphase transport of particulate cometary material as an explanation of the smooth terrains and repetitive outbursts on 9P/Tempel 1. *Icarus* **200**, 280–291 (2009).
- Belton, M. J. S. *et al.* The internal structure of Jupiter family cometary nuclei from Deep Impact observations: The “talps” or “layered pile” model. *Icarus* **187**, 332–344 (2007).
- Lasue, J., Botet, R., Levasseur-Regourd, A. C. & Hadamcik, E. Cometary nuclei internal structure from early aggregation simulations. *Icarus* **203**, 599–609 (2009).
- Sierks, H. *et al.* On the nucleus structure and activity of comet 67P/Churyumov-Gerasimenko. *Science* **347**, aaa1044 (2015).
- Keller, H. U. *et al.* OSIRIS—the scientific camera system onboard Rosetta. *Space Sci. Rev.* **128**, 433–506 (2007).
- Thomas, N. *et al.* The morphological diversity of comet 67P/Churyumov-Gerasimenko. *Science* **347**, aaa0440 (2015).
- Thomas, P. C. *et al.* The shape, topography, and geology of Tempel 1 from Deep Impact observations. *Icarus* **187**, 4–15 (2007).
- Thomas, P. C. *et al.* The nucleus of Comet 9P/Tempel 1: shape and geology from two flybys. *Icarus* **222**, 453–466 (2013).
- Brownlee, D. E. *et al.* Surface of young Jupiter family comet 81P/Wild 2: view from the Stardust spacecraft. *Science* **304**, 1764–1769 (2004).
- Britt, D. T. *et al.* The morphology and surface processes of comet 19/P Borrelly. *Icarus* **167**, 45–53 (2004).
- A'Hearn, M. F. *et al.* Deep Impact: excavating comet Tempel 1. *Science* **310**, 258–264 (2005).
- A'Hearn, M. F. *et al.* EPOXI at comet Hartley 2. *Science* **332**, 1396–1400 (2011).
- Wahlberg Jansson, K. & Johansen, A. Astrophysics formation of pebble-pile planetesimals. *Astron. Astrophys.* **570**, A47 (2014).
- Jutzi, M. & Asphaug, E. The shape and structure of cometary nuclei as a result of low-velocity accretion. *Science* **348**, 1355–1358 (2015).
- Capaccioni, F. *et al.* The organic-rich surface of comet 67P/Churyumov-Gerasimenko as seen by VIRTIS/Rosetta. *Science* **347**, aaa0628 (2014).

Supplementary Information is available in the online version of the paper.

Acknowledgements OSIRIS was built by a consortium of the Max-Planck-Institut für Sonnensystemforschung (in Göttingen, Germany), CISAS-University of Padova (Italy), the Laboratoire d'Astrophysique de Marseille (France), the Instituto de Astrofísica de Andalucía, CSIC (Granada, Spain), the Research and Scientific Support Department of the European Space Agency (Noordwijk, The Netherlands), the Instituto Nacional de Técnica Aeroespacial (Madrid, Spain), the Universidad Politécnica de Madrid (Spain), the Department of Physics and Astronomy of Uppsala University (Sweden), and the Institut für Datentechnik und Kommunikationsnetze der Technischen Universität Braunschweig (Germany). The support of the national funding agencies of Germany (DLR), France (CNES), Italy (ASI), Spain (MEC), Sweden (SNSB), and the ESA Technical Directorate is gratefully acknowledged. We thank the ESA teams at ESAC, ESOC and ESTEC for their work in support of the Rosetta mission.

Author Contributions M.M. led and designed the study, identified and mapped most of the strata, selected the areas for retrieving 3D best-fitting planes, performed the geological sections, made the overall geological interpretation and wrote most of the text; E.S. carried out the 3D reconstruction of strata attitudes and wrote part of the main text and Methods; F.M. obtained the gravity-field vectors, wrote part of the Methods and contributed to data interpretation; G.C. contributed to designing the study and data interpretation; L.G. performed the detailed geomorphological analysis of the Hatmehit region; M.P. contributed to the geomorphological analysis of the Seth region and the landing-site candidate A; L.J. was responsible for the stereo-photoclinometric model and the 3D reconstruction of the two lobes as independent objects; G.N. and S.L. substantially contributed to data interpretation and defining the related implications; F.P. and F.S. were responsible for the stereo-photogrammetric shape model; M.R.E.-M. defined the physiographic regions of the comet. H.S., C.B., P.L., R.R., D.K. and H.R. are the lead scientists of the OSIRIS project. The other authors are all co-investigators who built and ran this instrument and made the observations possible, and associates and assistants who participated in the study.

Author Information All data presented in this paper will be delivered to ESA's Planetary Science Archive (<http://www.rssd.esa.int/index.php?project=PSA&page=rosetta>) and NASA's Planetary Data System (<https://pds.nasa.gov/>) in accordance with the schedule established by the Rosetta project. Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to M.M. (matteo.massironi@unipd.it).

METHODS

Data. For the observational work and geological mapping we used both OSIRIS NAC and WAC images. The NAC images used for this work were acquired soon after the orbit insertion on 6 August 2014 up to 17 March 2015 and have a spatial scale ranging from about 2 m per pixel to 0.1 m per pixel depending on the cometocentric distance of acquisition (0.1 m per pixel for the 14 February 2015 close fly-by); WAC data were considered only for distances less than 45 km and have a spatial scale ranging from about 4.5 m per pixel to 2.5 m per pixel (Extended Data Table 1). Lower-resolution NAC images and all WAC images provided synoptic views of regional sectors of the whole nucleus and were useful for the broad mapping of strata and terraces. High-resolution NAC data were instead used to unravel details of the stratigraphic sequences and their presence within pits and along steep walls in several locations at the nucleus surface.

Surface strata description. (See also Supplementary Information section 'Concept of strata applied to the 67P/C-G comet nucleus'). The deep depression (the 'neck') separating the larger lobe (the 'main body') from the smaller one (the 'head') allows for detailed views of both sides, providing insight into the inner structure of the two lobes. For this reason we describe the main stratified structure of the comet starting from the neck and going towards the main body and the head. Our observations refer to the regional physiographical distinctions defined by the OSIRIS team^{7,17}. These morphological regions are probably the results of uneven fracturing, weathering and sublimation processes on the comet. The differential action of these processes on the primary structure most probably shaped the regions now visible on the comet nucleus¹⁷.

The main body has a size of 4.1 km × 3.3 km × 1.8 km along three perpendicular directions⁵ and shows several morphologies indicating inner stratification. The most remarkable ones are the terraces and cliffs of the Seth region facing the neck region (Hapi) (Fig. 1a).

The Seth region includes terraces at different scales; these are always aligned with other terraces and lineaments, indicating a continuous inner stratification (Fig. 1a and Extended Data Fig. 1a–c). An impressive example is given by the terrace at the candidate landing site labelled A (a plain 400 m long and 600 m wide delimited by two steep walls) (Extended Data Fig. 1d). The upper wall (400 m high) shows smaller-scale aligned terraces joined by linear discontinuities that are stratigraphically traceable for hundreds of metres. These small terraces have orientations similar to the main terrace and coherent with the lineaments of the lower wall (300 m high). On one side of the main terrace there is a pit 220 m wide and 185 m deep, characterized by strata as thin as 3 m that evidently represent the inner skeleton of the comet nucleus up to at least some hundreds of metres deep (Fig. 1b). Again, the orientation of this stratified sequence is compatible with the one visible in the main terrace and the walls. Hence, the average thickness enclosed between terraces is not representative of the thickness of the elemental strata, which is much thinner, a few tens of centimetres to a few metres (compatible with the particle sizes of the accreting materials¹⁴).

The Seth terraces shown by the OSIRIS images appear to gradually change their orientation, that is, they become more inclined and 'dip' towards the neck region (Hapi), and are less inclined towards the Ash region. This general trend is illustrated by the lateral boundary separating Seth and Ash from Anubis and Atum (Fig. 1c). This boundary consists of a continuous scarp whose maximum elevation apparently constrains the stratigraphic thickness of the Seth region to 150 m. However, the presence of terraces in Anubis underneath the smooth deposits (Fig. 1c) suggests that the stratification probably reaches greater depths from the surface.

The Ash region, which connects Seth (facing the head) to Imhotep (on the opposite site), is characterized by several cuesta-like and mesa-like features (Fig. 1d and Extended Data Fig. 2a–c). On Earth, cuestas are formed by differential erosion of a stratified sequence of differently erodible strata inclined less than 30° with respect to the surrounding average topography¹⁸. The orientation of the strata is coincident with the dip and the dip direction (that is, the direction of maximum inclination) of the gentle slope. Mesas originate via the same process but are characterized by horizontal terraces. Similarly, the cuesta-like and mesa-like morphologies on the comet nucleus could have been generated by a differential erosional sublimation of a stratified sequence of strata with different volatile contents.

Across the Ash region, the strata gradually change their dip direction, assuming an inclination opposite to that of the Seth terraces and slightly dipping towards the Imhotep depression. The low dip angles of Ash stratification favour the formation of sequences of aligned flat terraces that can be followed for hundreds of metres to kilometres in mesa-like morphologies (Fig. 1d and Extended Data Fig. 2a–c).

At the opposite site of Anubis, Seth and Ash grade into the Babi region⁷. The gradual passage between these provinces means that no major change and break occur in the inner structure, suggesting it is constituted by a continuous inner stratigraphy. Indeed, Babi is characterized by terraces and strata at all scales (Extended Data Fig. 3).

The Imhotep depression is entirely surrounded by scarps with evidence of stratification (Fig. 1d and Extended Data Fig. 2a and b), suggesting a stratified structure for all the confining provinces. Hence, Kephry and Apis also show evidence of a pervasive inner stratification. The elevation of the Imhotep bounding scarps is around 200 m, and this would roughly define the minimum thickness of the stratified structure in the surrounding regions. Within Imhotep itself there are several orders of staircase terraces underneath smooth materials physically connected with the strata at Imhotep boundaries (Fig. 1e and Extended Data Fig. 2a, b and d). Strata within Imhotep must be considered as a fundamental structure.

Surface-emplaced geological units on comets could be sequences of air-fall deposits of lag materials¹ and successive flows of fluidized materials². Strata of both air-fall deposits and fluidized flows should display smooth and gradual terminations, and not sharp boundaries as do the steps (terrace margins), shown in metre-to-centimetre scale in Fig. 1e and Extended Data Figs 2d and e. On the other hand, later erosive processes might form angular boundaries on previous smooth terminations. However, flow features have not been detected at any scale on OSIRIS images of the Imhotep area and this observation definitively rules out the Belton and Melosh mechanism² of emplacement in this region. Moreover, the terrace margins within the Imhotep depression are connected by flat planes covering several square kilometres; any air-fall deposit would instead mantle the very irregular pre-existing topography, making the formation of such continuous flat areas very unlikely. All of this evidence proves that air-fall deposits, represented by the loose material on top of the stratified sequence, cannot explain the staircase succession of terraces in Imhotep. These terraces clearly connect the terrace margins with the strata of the bedrock walls bounding the Imhotep region, and even the lateral prosecution of the terraces in the smooth terrain are connected to the bedrock stratification (Fig. 1e). In Extended Data Fig. 2d the strata forming the largest roundish plateau at the Imhotep boundary must be exhumed, since they dip underneath the loose mantling material. All these observations mean that the stratified inner structure is thicker than is expected from looking at the average elevation of the scarp bounding the Imhotep depression.

To summarize, all regions of the main body are characterized by strata and terraces with the exception of Aten, which is an elongated depression covered by debris. Several different views (Fig. 1c, d and Extended Data Fig. 2a and b) suggest that the comet is at least partially wrapped in strata.

The most spectacular section across the minor lobe (2.6 km × 2.3 km × 1.8 km in size⁵) is Hathor, a wall 900 m high and 1,500 m wide displaying about one third of the head's inner structure⁷ and facing the Hapi region. Hathor cliff is made up of a highly pervasive stratification demonstrated by lineaments with a regular spacing and intimately associated with small terraces (Fig. 2a). Strata are crosscut by a perpendicular set of lineaments that can be interpreted as fracture systems since they are not regularly spaced, are not associated to any terrace, and are normally continuous and not interrupted by strata.

At the sharp topographic boundary between Hathor and Ma'at no abrupt change of strata orientations is recorded (Fig. 2a). In contrast, the same stratification seems to be involved in a change in orientation on both sides of the main wall of Hathor, where they are steeply inclined (Fig. 2a). However, the small terraces within Hathor are slightly dipping towards Hapi as well and, at different views, the lateral variation on strata attitudes between Hathor and the confining Anuket region appears to be continuous. The alcove carved at the contact between Hathor and Anuket^{7,17} marks the change in the strata orientation, providing a good lateral section of what comprises the Anuket inner structure (Fig. 2b). As usual for any topographic surface parallel to strata, Anuket (as well as Apis and Kephry in the main body; Extended Data Fig. 2b) displays only some, although meaningful, examples of terraces (Fig. 2b and c). Thomas *et al.*⁷ have pointed out how Anuket materials seem to extend underneath the Ma'at region (Fig. 2c). This explains the continuous change in orientation of strata between Anuket and Hathor/Ma'at, and inherently implies that the stratification visible on the Hathor 900-m-high wall is actually composed, in its upper portion, of Anuket and then Ma'at materials. At the same time, this seems to reinforce the observation that Hathor is representative of a wide part of the head's inner structure that has been exposed because of the dismantling (by erosion or detachment) of a wide portion of the head whose remnant is the distinctive Anuket region that we see today.

Like Ash in the main body, Ma'at, Maflet and Bastet are also characterized by cuesta- and mesa-like morphologies, often affected by gravitational niches and distributed in staircase patterns (Extended Data Fig. 4a). The orientation of the related terraces varies very gradually across these regions and no abrupt changes are recorded. The Hatmehit depression is carved from these regions and its bounding scarps give good views into the subsurface structures of the confining regions. In particular, in OSIRIS context images centripetal patterns are apparent, mainly at the boundary between Hatmehit and Ma'at (Extended Data Fig. 4a). However, at a closer view these features appear to be the result of composite morphologies resulting from the interplay of inclined terraces caused by dipping

strata and fractures (Extended Data Fig. 4b and c). This further indicates that Ma'at, Maflet and Bastet contain a stratified sequence at least as thick as the depth of the Hatmehit depression (around 150 m). However, the prominent step revealed within Hatmehit, which does not extend outside the depression (Extended Data Fig. 4a), appears to be the expression of a buried terrace rather than of a later deformation (Extended Data Fig. 4d and e). This would imply a stratified sequence thicker than the maximum elevation of the scarps over the depression.

Finally, the Nut and Serqet regions can be considered from a structural point of view as a unique system. In fact, the Serqet scarp facing Nut displays strata emphasized by small terraces whose orientations are parallel to the overall trend of the Nut flat surface representative of a large strata plane itself (Fig. 2d).

In summary, all regions of the minor lobe except Hatmehit, which shows a buried terrace, appear to be characterized by stratified sequences. Moreover, the structural relationship between Hathor, Ma'at and Anuket as well as the orientation of terraces bounding Hatmehit and the cuesta-like morphologies of the other regions suggest that a stratified envelope characterizes the head as well.

Best-fitting planes. All the procedures of plane extraction and 3D geological reconstructions of strata attitudes were implemented using MeshLab (<http://meshlab.sourceforge.net/>) and Matlab software on a shape model obtained through stereo-photoclinometry and covering 75% of the entire surface. The shape model is characterized by 3.6 million facets at 6 m sampling. The best-fitting planes were extracted on the bases of the shape nodes of the terraces. The number of points used for best-fitting is 420 on average. Extended Data Figure 5 shows the distribution of the ratios between the standard deviations of the 3D coordinates in the local reference system (where the z axis is normal to the surface). The mean ratio is 0.0574, which corresponds to a maximum error on the dip angle and the dip direction of the best-fitting planes lower than 2° .

Geological sections. The stereo-photoclinometric shape model has been the source for retrieving cross-sectional topographic profiles parallel (longitudinal) and perpendicular (transversal) to the major axis of the comet nucleus. On the derived topographic profiles we projected all the terrace-related planes (that is, best-fitting planes) whose centres were at a distance not exceeding 500 m from the plane of the section. (A similar projection of surface measurements into specific topographic profiles was used for the geological sections of the 67-km-long Brenner Basis Tunnel^{19,20}, which will cross through the European Alps; <http://www.bbt-se.com/en/home/>.) The projected best-fitting planes are red in the top panels of Fig. 4 and the left panels of Extended Data Figs 7–9 and blue/green in the related cross-sectional slices.

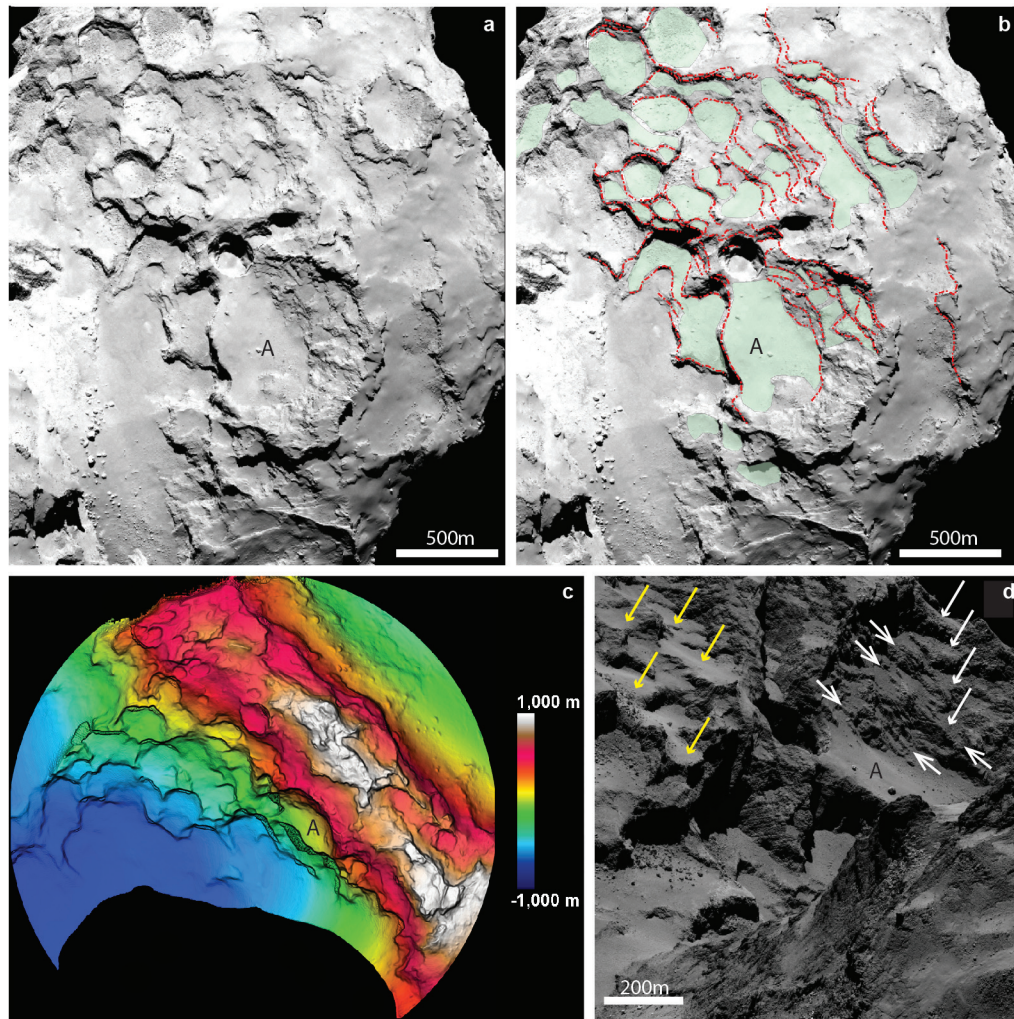
Following consolidated stratigraphic procedures²¹, we have considered the strata to be parallel to one other, laterally continuous, and with a constant thickness throughout the section. This means that the attitude of inner strata is constrained by the orientation of the strata above at the surface. Hence vectors perpendicular to the strata planes at the surface interpolated from one terrace to another define the field lines (yellow in the cross-sectional slices), which constrain the strata orientations underneath the surface (black dashed lines). By drawing the geological sections starting from the best-fitting planes transecting the

topographic profile and following the constraints given by the derived field lines (that is, strata must be perpendicular to the yellow field lines), we were able to partially reconstruct the geological structure of the interior of the nucleus up to depths of several hundred metres (Fig. 4 and Extended Data Figs 7–9). In sum, all the strata drawn in the sections originate from one or more measurements taken at the surface, are parallel to each other and laterally continuous (as they are in the OSIRIS images), maintain the same thickness throughout the section and have their attitude constrained by the orientation of the strata above. The strata identified in the longitudinal and transversal sections are not necessarily the same since they cut through different topographies.

Figure 4 and Extended Data Fig. 9 show all the strata intersected by the displayed sections, whereas Extended Data Figs 7 and 8 show simplified longitudinal sections, to clarify the procedure we used. Indeed, on one hand, within the longitudinal sections of Extended Data Figs 7 and 8, we have reported only the extension of the strata measured at the surface of the same section (blue lines transecting the topographic profiles). On the other hand, in each transversal section of Extended Data Fig. 9 we have considered all the strata, including the ones derived only from the longitudinal sections and intersecting the transversal traces at depth (Extended Data Fig. 9). In summary, almost the whole interpretation has been directly constrained by the orientation of the best-fitting planes and related field lines. Geo-structural interpretation indirectly controlled by the attitudes of contiguous strata and mere observations of the surface, which is a common practice in drawing classical geological sections of the Earth and other planets of the Solar System, has been applied in rare locations and is identified by grey lines in Fig. 4 and Extended Data Figs 7–9.

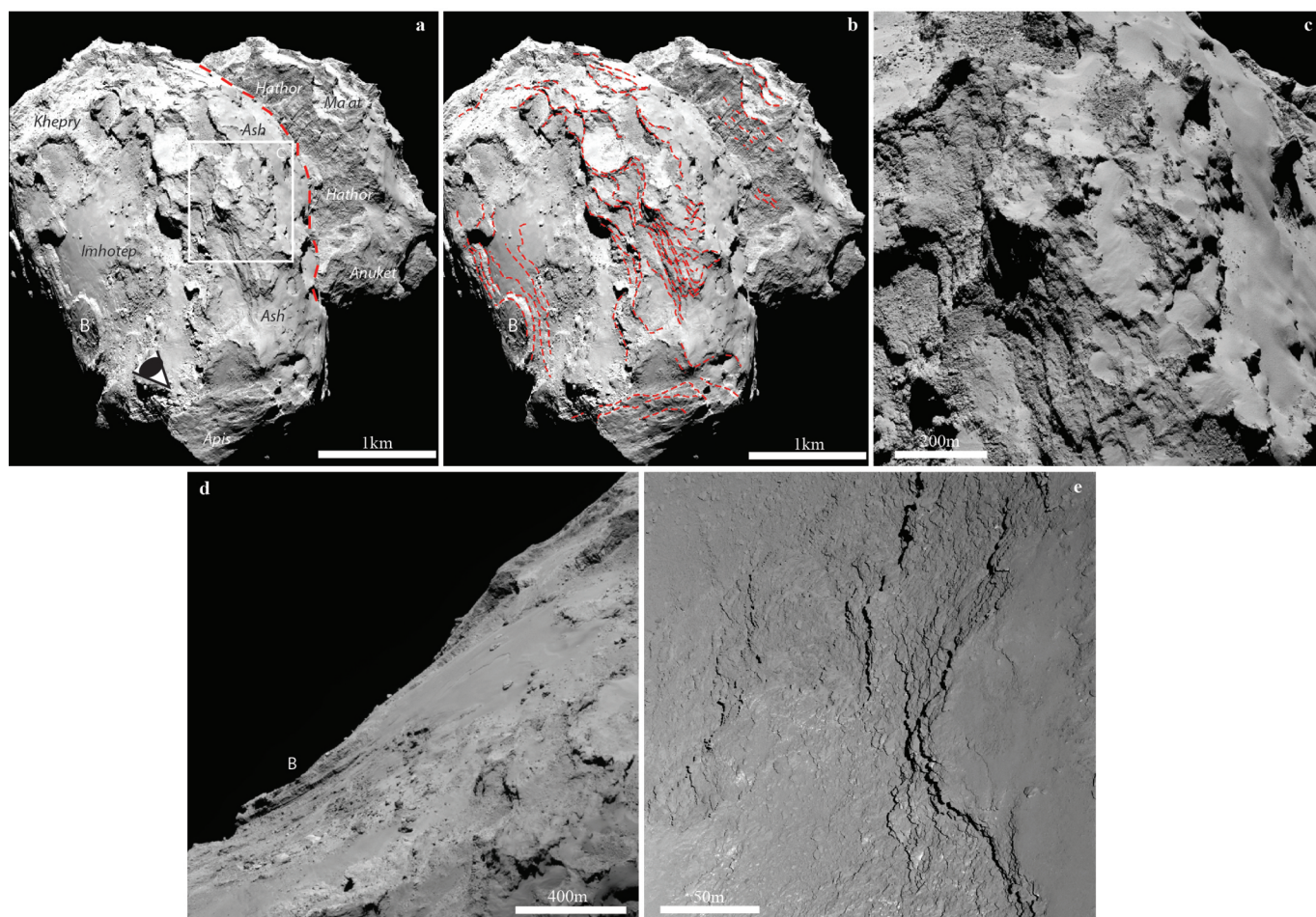
Gravity-field vectors. The local gravity-field vectors of the entire comet nucleus and of the two separated lobes have been computed by modelling the bodies as polyhedrons with triangular facets. The gravity potential is computed for each shape model, following the theory of Werner and Scheeres²², adding up the contribution from all the facets of the polyhedron. The centrifugal force due to the comet rotation is included in the final computation of the force field.

17. El-Maarry, M. R. *et al.* Regional surface morphology of comet 67P/Churyumov-Gerasimenko from Rosetta/OSIRIS images. *Astron. Astrophys.* <http://dx.doi.org/10.1051/0004-6361/201525723> (2015).
18. Strahler, A. N. *Physical Geography* (John Wiley & Sons, 1975).
19. Bistacchi, A. *et al.* 3D fold and fault reconstruction with an uncertainty model: an example from an Alpine tunnel case study. *Comput. Geosci.* **34**, 351–372 (2008).
20. Schiavo, A. *et al.* Geology of the Brenner pass-Fortezza transect, Italian eastern Alps. *J. Maps* **11**, 201–215 (2015).
21. Groshong, R. H. *3-D Structural Geology. A Practical Guide to Quantitative Surface and Subsurface Map Interpretation* (Springer, 2008).
22. Werner, R. A. & Scheeres, D. J. Exterior gravitation of a polyhedron derived and compared with harmonic and mascon gravitation representations of asteroid 4769 Castalia. *Celestial Mech. Dyn. Astron.* **65**, 313–344 (1997).
23. Preusker, F. *et al.* Shape model, reference system definition, and cartographic mapping standards for comet 67P/Churyumov-Gerasimenko. Stereo-photogrammetric analysis of Rosetta/OSIRIS image data. *Astron. Astrophys.* <http://dx.doi.org/10.1051/0004-6361/201526349> (2015).



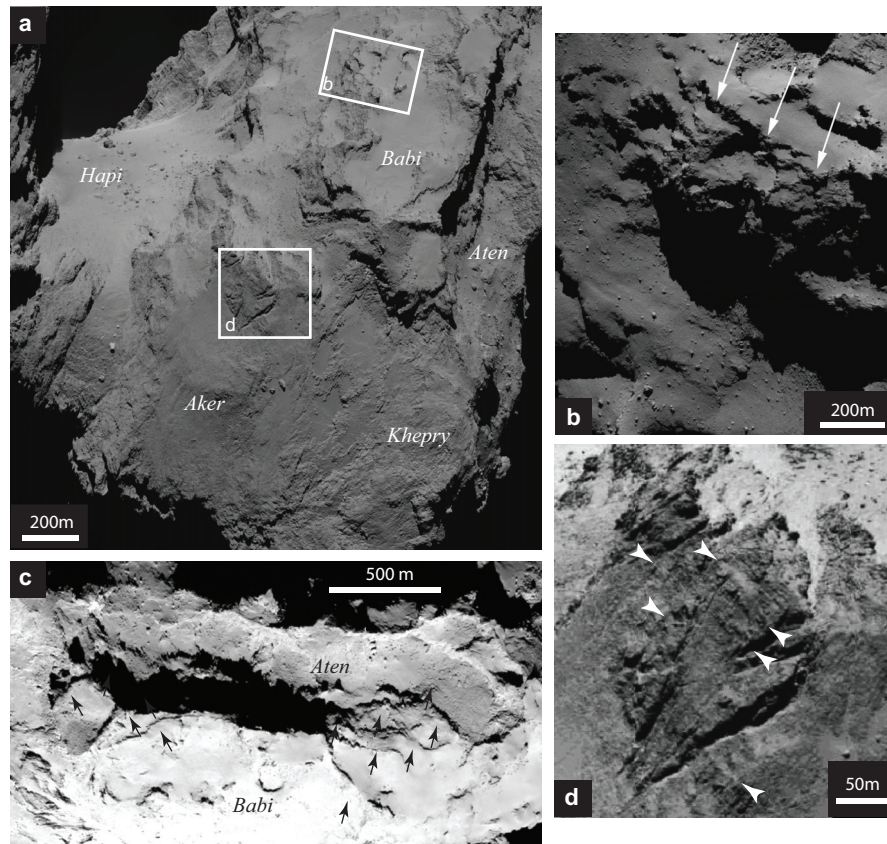
Extended Data Figure 1 | Seth region. **a, b,** General view (**a**) and interpreted image (**b**) with the main terraces (pale green), and strata margins (red dashed lines) of the Seth region. **c,** Stereographic projection (north-polar aspect) of the main lobe, displaying heights derived from a stereo-photogrammetric shape model²³. Colour scale indicates height above a sphere of radius 1.5 km centred

in the main lobe. Terraces are aligned along different levels of the same elevation. **d,** NAC view of landing site candidate A (see location in panels **a–c**). Terraces in the upper wall of site A (small white arrows) together with parallel lineaments (large white arrows), define a continuous stratification. Yellow arrows indicate other terraces of the Seth region.



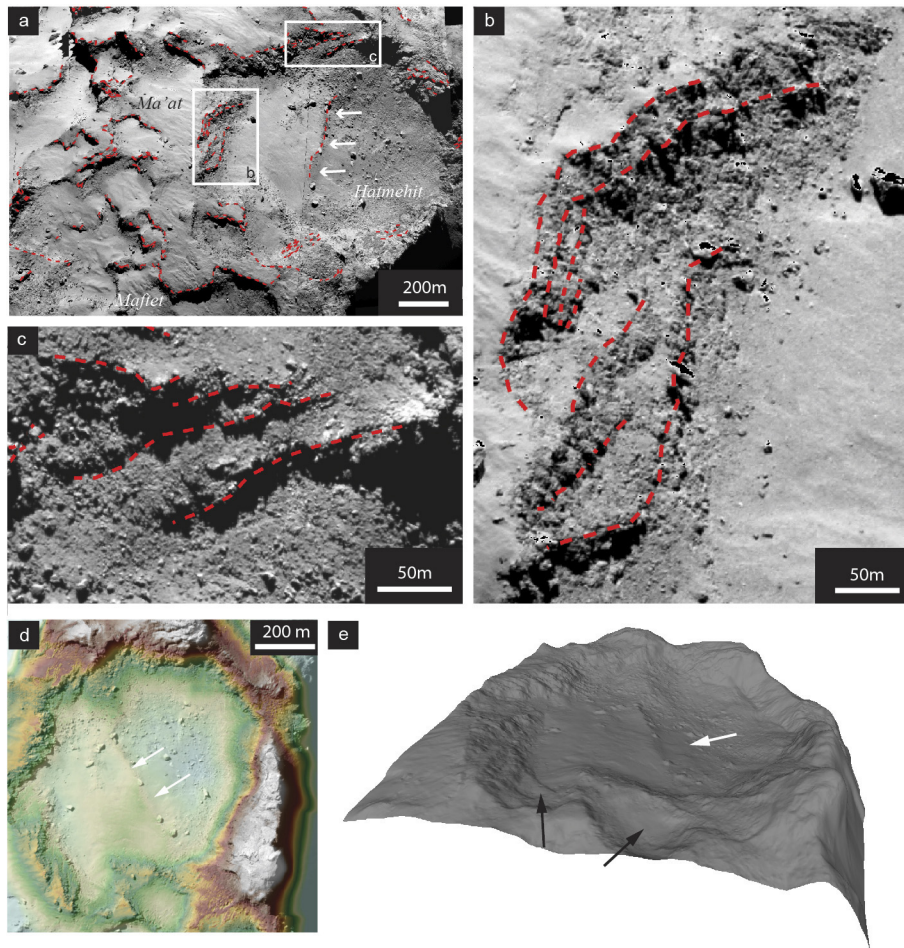
Extended Data Figure 2 | Imhotep regions and surroundings. **a**, NAC overview of the Ash, Apis, Khepry and Imhotep regions; in the background are regions of the head. The red dashed line separates the head from the main body. The white square is the location of **c**. The eye symbol shows the view of **d**. **b**, Map of the terrace, cuesta and mesa margins (dashed red lines) showing

the general attitudes of strata in the main body. **c**, Close view of the stratification within Ash, where thinner strata are only a few metres thick. **d**, Imhotep terrace margins and inner stratification. Strata of site B dip underneath smooth deposits, showing their exhumed nature. **e**, Metre-to-centimetre scale strata heads within Imhotep.



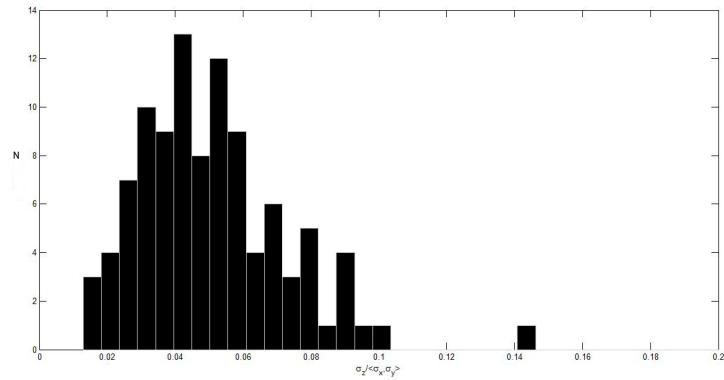
Extended Data Figure 3 | The Babi and Aten regions. **a**, NAC overview of Babi region and surroundings. White squares are the locations of **b** and **d**. **b**, Detailed NAC view of Babi stratification: white arrows indicate small

terraces and strata. **c**, The Babi–Aten boundary is underlined by staircase terraces, indicated by black arrows. **d**, Detailed view of stratification in Aker: white arrows show strata cut by fractures.



Extended Data Figure 4 | The Hatmehit region. **a**, NAC mosaic of the Hatmehit region and surroundings. The red dashed lines show the cuesta-like margins in Ma'at and Maftet. White arrows indicate the terrace step within Hatmehit. White squares are the locations of **b** and **c**. The apparently centripetal morphologies all around Hatmehit are due to stratified scarps,

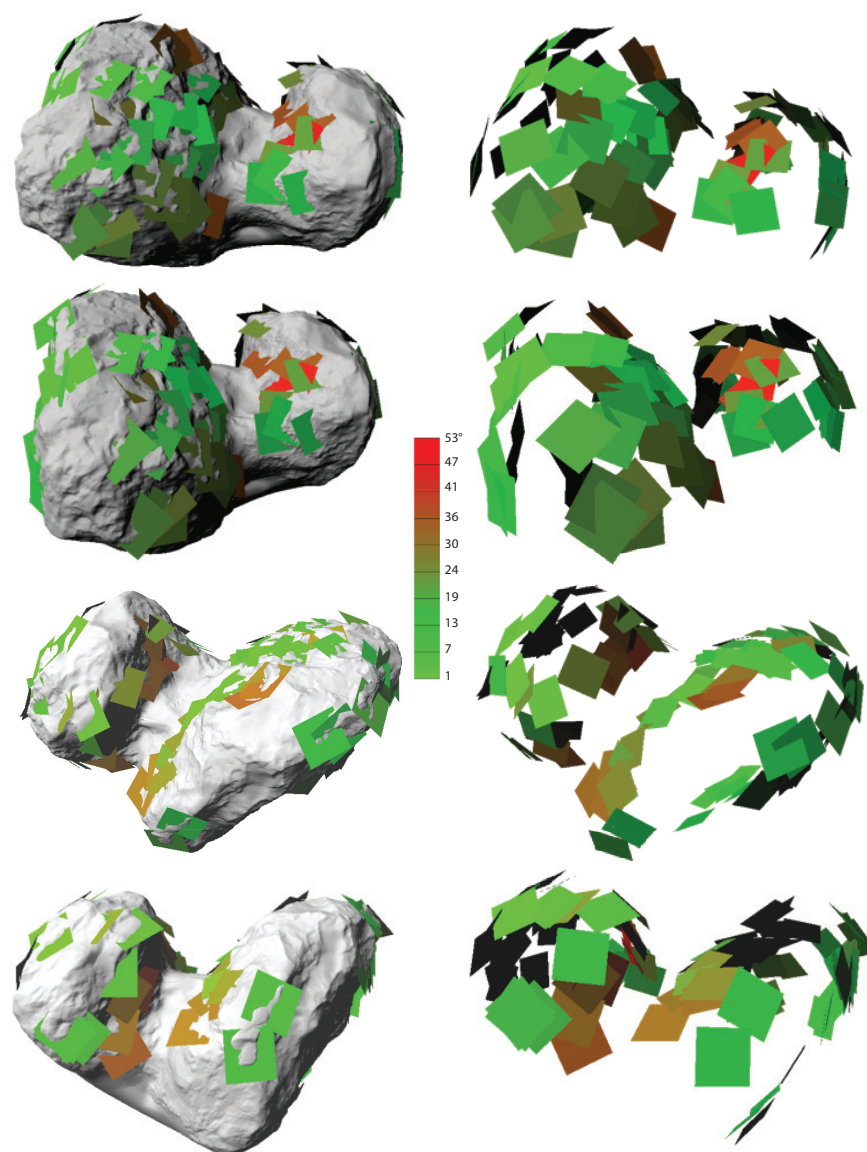
highlighted by the red dashed lines in **b** and **c**. The terrace step within Hatmehit (white arrows in **d** and **e**) and the terraces at its margin (black arrows in **e**) are visible in the stereo-photoclinometry Digital Terrain Model (**d**) and the related perspective view (**e**).



Extended Data Figure 5 | Standard deviations of best-fitting planes.

Distribution of the ratios between the standard deviations of the 3D coordinates in the local reference system (where the z axis is normal to the best-fitting plane). N is the number of best fitting planes; σ_z is the standard deviation along

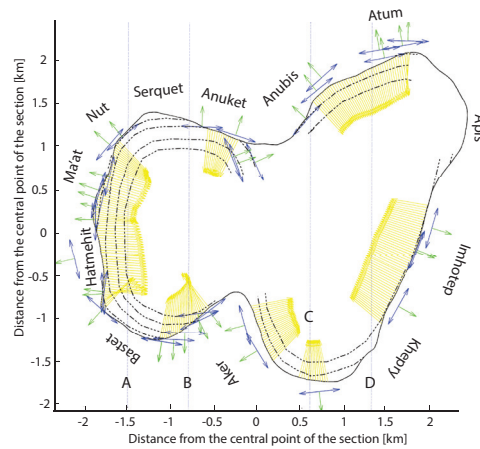
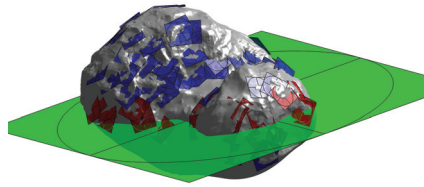
the z axis of the points used to retrieve a given plane; σ_x and σ_y are the standard deviations along the other two perpendicular axes. The mean ratio is 0.0574, which corresponds to a maximum error on the dip angle and dip direction of the best-fitting planes lower than 2° .



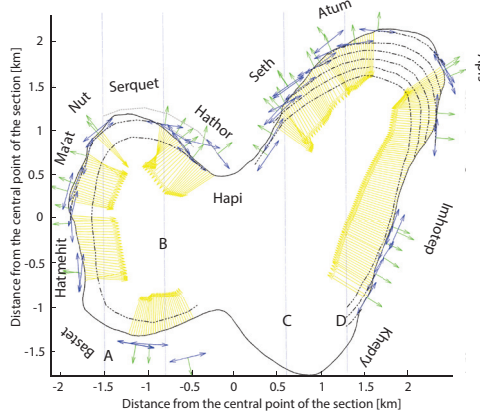
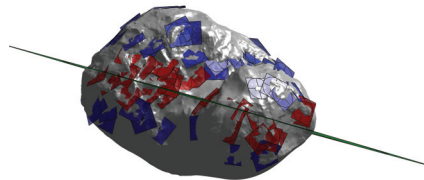
Extended Data Figure 6 | 3D views of the nucleus and strata of comet 67P/Churyumov-Gerasimenko. On the left side are four different 3D views of the nucleus of comet 67P/Churyumov-Gerasimenko with best-fitting planes derived from the stratification. On the right side are the best-fitting planes alone. Each plane indicates the orientation of strata at that specific location

(corresponding to the centre of the drawn plane) on the comet nucleus. Note how the two lobes show independent envelopes. The colour scale indicates the angular deviation between the plane vector and the local gravity vector (calculated for the whole body, assuming uniform density).

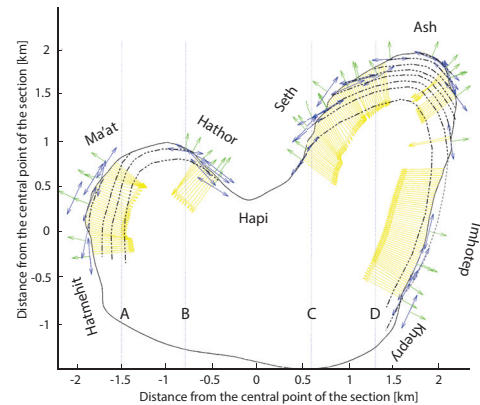
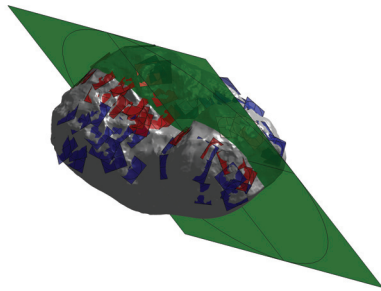
1



2

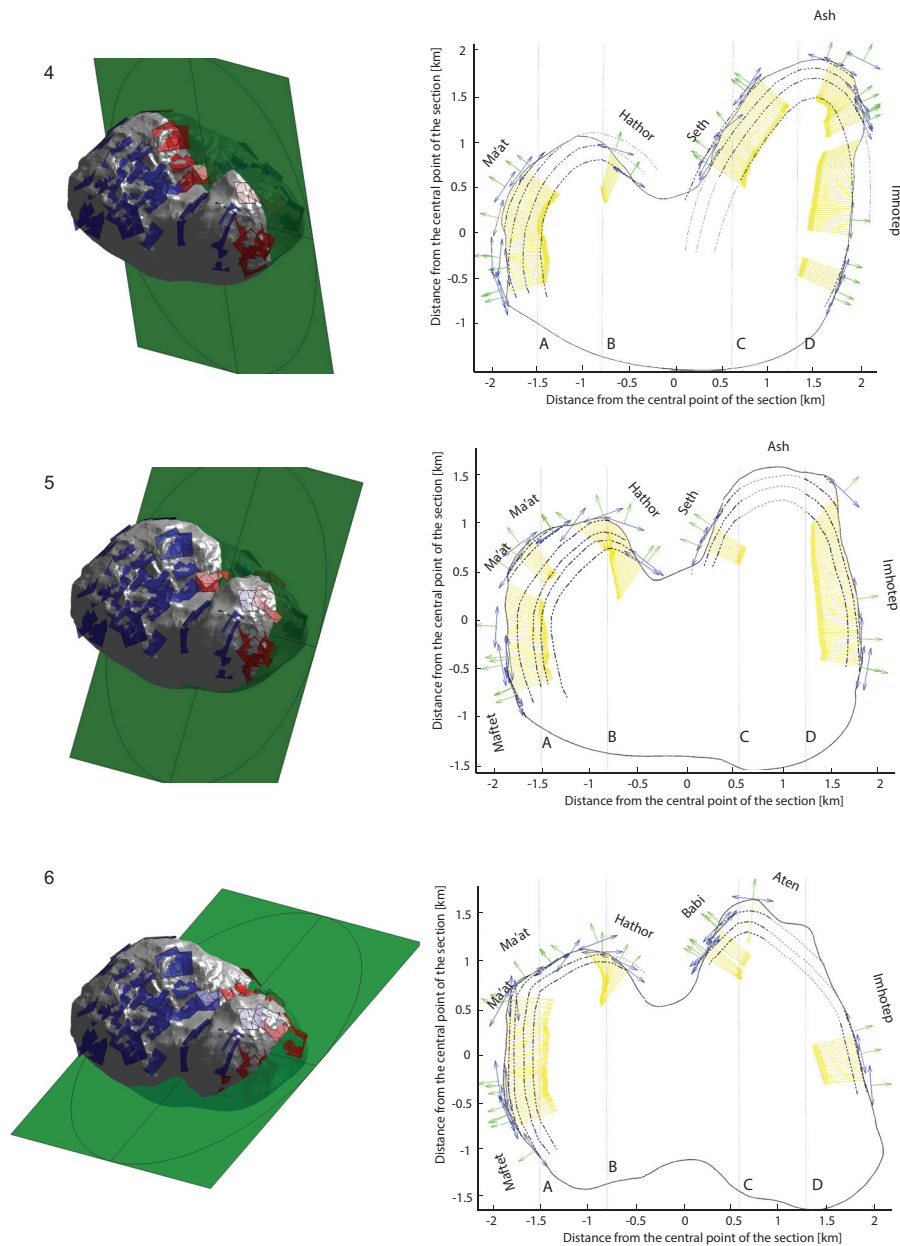


3



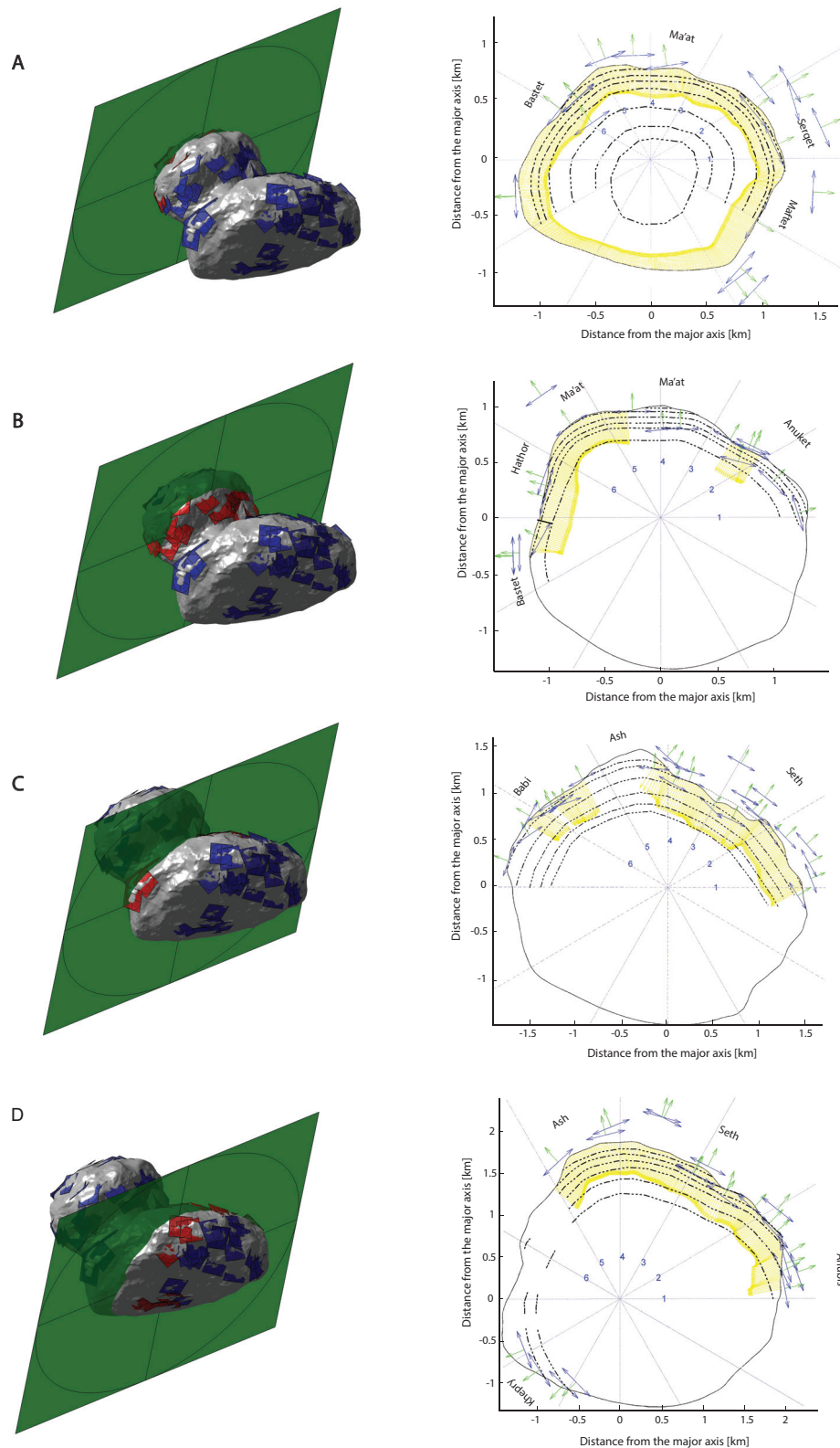
Extended Data Figure 7 | Eastern sector cross-sectional slices through comet 67P/Churyumov-Gerasimenko along its major axis (longitudinal sections). On the left are perspective views of the comet nucleus showing the different cross-sectional slices. In red are shown the best-fitting planes projected into the geological section; in blue all the other best-fitting planes. On the right are geological sections. In blue are shown the best-fitting planes; green

arrows are vectors perpendicular to each best-fitting plane; yellow lines show the field of lines used for drawing strata within the comet nucleus (strata are perpendicular to the yellow lines); dashed black lines are the strata departing from measurements at the surface of each longitudinal section. The traces of transversal sections A, B, C and D and the locations of some regions are also reported.



Extended Data Figure 8 | Western sector cross-sectional slices through the comet along its major axis (longitudinal sections). On the left are perspective views of the comet nucleus showing the different cross-sectional slices. In red are shown the best-fitting planes projected into the geological section, in blue is shown all the other best-fitting planes. On the right are geological sections. In blue are shown the best-fitting planes; green arrows are vectors

perpendicular to each best-fitting plane; yellow lines show the field of lines used for drawing strata within the comet nucleus (strata are perpendicular to the yellow lines); dashed black lines are strata departing from measurements at the surface of each longitudinal section; dashed grey lines are inferred strata. The traces of transversal sections A, B, C and D and the locations of some regions are also reported.



Extended Data Figure 9 | Cross-sectional slices through the comet perpendicular to the major axis (transversal sections). On the left are perspective views of the comet nucleus showing the different cross-sectional slices. In red are shown the best-fitting planes projected into the geological section; in blue are shown all the other best-fitting planes. On the right are geological sections. In blue are shown the best-fitting planes; green arrows are vectors perpendicular to each best-fitting plane; yellow lines show the field of

lines used for drawing strata within the comet nucleus (strata are perpendicular to the yellow lines); dashed black lines are strata departing from measurements at the surface of each section and from the intercept of strata of the longitudinal sections; dashed grey lines are inferred strata. The traces of longitudinal sections 1, 2, 3, 4, 5 and 6 and the locations of some regions are also reported.

Extended Data Table 1 | ID numbers for the OSIRIS images

Figure	OSIRIS IMAGES
1a	NAC_2014-08-06T02.20.01_F28
1b	NAC_2014-09-22T01.48.46_F41
1c	NAC_2014-09-05T02.47.24_F71
1d	WAC_2014-09-18T05.23.28_F17
1e	NAC_2015-03-10T20.05.49_F22
2a	NAC_2014-08-07T20.37.34_F22
2b	NAC_2014-09-13T08.19.08_F22
2c	WAC_2014-09-12T18.51.33_F17
2d	NAC_2014-09-19T12.41.05_F41
Extended Data	OSIRIS IMAGES
1a,b	NAC_2014-08-08T03.37.34_F22
1d	NAC_2014-09-19T13.16.13_F22
2a,b	NAC_2014-08-16T14.59.14_F22
2c	NAC_2014-09-20T19.01.50_F22
2d	NAC_2015-03-17T03.37.46_F41
2e	NAC_2015-02-14T12.36.35_F88
3a	NAC_2014-08-25T12.42.54_F22
3b	NAC_2014-09-13T01.02.49_F22
3c	NAC_2014-08-16T17.59.14_F22
3d	NAC_2014-09-05T09.35.5_F22
4a,b,c	Mosaic: NAC_2014-09-14T17.47.55_F22, NAC_2014-09-14T17.32.15_F41, NAC_2014-09- 14T17.31.24_F22, NAC_2014-09-15T05.43.04_F41, NAC_2014-09-14T05.05.28.432_F22

The first three letters indicate the instrument used to acquire the image (NAC or WAC); the following digits are the time (in Coordinated Universal Time, UTC) of imaging (year-month-day, then hour-minute-seconds); the last two numbers correspond to the filters.

Observation of the competitive double-gamma nuclear decay

C. Walz¹, H. Scheit¹, N. Pietralla¹, T. Aumann¹, R. Lefol^{1,2} & V. Yu. Ponomarev¹

The double-gamma ($\gamma\gamma$)-decay of a quantum system in an excited state is a fundamental second-order process of quantum electrodynamics. In contrast to the well-known single-gamma (γ)-decay, the $\gamma\gamma$ -decay is characterized by the simultaneous emission of two γ quanta, each with a continuous energy spectrum. In nuclear physics, this exotic decay mode has only been observed for transitions between states with spin-parity quantum numbers $J^\pi = 0^+$ (refs 1–3). Single-gamma decays—the main experimental obstacle to observing the $\gamma\gamma$ -decay—are strictly forbidden for these $0^+ \rightarrow 0^+$ transitions. Here we report the observation of the $\gamma\gamma$ -decay of an excited nuclear state ($J^\pi = 11/2^-$) that is directly competing with an allowed γ -decay (to ground state $J^\pi = 3/2^+$). The branching ratio of the competitive $\gamma\gamma$ -decay of the $11/2^-$ isomer of ^{137}Ba to the ground state relative to its single γ -decay was determined to be $(2.05 \pm 0.37) \times 10^{-6}$. From the measured angular correlation and the shape of the energy spectra of the individual γ -rays, the contributing combinations of multipolarities of the γ radiation were determined. Transition matrix elements calculated using the quasiparticle-phonon model reproduce our measurements well. The $\gamma\gamma$ -decay rate gives access to so far unexplored important nuclear structure information, such as the generalized (off-diagonal) nuclear electric polarizabilities and magnetic susceptibilities³.

Two-photon processes, that is, excitation and decay of a quantum state simultaneously involving two photons whose energy sum $E_1 + E_2$ matches the transition energy E_0 , have been studied intensely in the past two decades in atomic physics^{4,5} and are now routinely applied in spectroscopic studies. These processes were first discussed theoretically in 1929^{6,7}, and the same approach predicted the existence of the rare double- β -decay of atomic nuclei as the analogous second-order process for the electroweak interaction⁸. In contrast to the situation in atomic physics, data on the $\gamma\gamma$ -decay of atomic nuclei are very sparse. Up to now, $\gamma\gamma$ -decays of nuclear states have only been observed for the special cases of the first excited states of the ‘doubly-magic’ nuclei ^{16}O , ^{40}Ca and ^{90}Zr (refs 1–3, 9). The first excited states of these nuclei have the unusual property of having spin-parity quantum numbers $J^\pi = 0^+$, which is also the case for their ground states. Hence, γ -decays are strictly forbidden for these transitions. For almost all other nuclei, the first excited states have spin quantum numbers larger than zero, making γ -decays possible. Despite several searches^{10,11}, the $\gamma\gamma$ -decay has never been observed in a situation where a single γ -ray transition is allowed, a situation we call the competitive $\gamma\gamma$ -decay, denoted by ‘ $\gamma\gamma/\gamma$ -decay’. Preliminary evidence for the $\gamma\gamma/\gamma$ -process obtained by other groups in parallel with our work has been announced elsewhere (ref. 12, and D. J. Millener and R. J. Sutter (personal communication); ref. 13, and C. J. Lister (personal communication)). In the present work, the experimental difficulties have been overcome and we report here firm observation of the $\gamma\gamma/\gamma$ -decay of a nuclear state.

The experimental challenge to the observation of the competitive $\gamma\gamma$ -decay arises from its decay rate, which is more than five orders of magnitude smaller than that of the allowed γ -decay. In an ensemble of nuclei undergoing γ -decay from a given excited state, there are a large

number of γ quanta being emitted, which cause two basic challenges, illustrated in Fig. 1. First, a γ quantum with transition energy E_0 can deposit part of its energy in one detector and then scatter into the second detector, depositing there its remaining energy. The sum of the energies registered in both detectors equals the transition energy, which is the signature of the $\gamma\gamma$ -decay. Second, two γ quanta, each with energy E_0 , can be emitted independently by two nuclei at almost the same time, and deposit a fraction of their energy in two detectors, such that the energy sum is again close to E_0 . Although these random coincidences can be subtracted from truly simultaneous events, they cause a substantial statistical uncertainty, thereby preventing the observation of the much rarer $\gamma\gamma/\gamma$ -decay unless they are suppressed properly.

In our experiment, the first excited $11/2^-$ state of ^{137}Ba was populated through β^- decays of ^{137}Cs nuclei, a well known γ -ray calibration standard. This isomer decays dominantly via a single- γ magnetic hexadecapole (M4) transition to the $3/2^+$ ground state, emitting a γ -ray with an energy of 661.66 keV (Fig. 2a inset). The emitted γ -rays were detected using five large-volume $\text{LaBr}_3:\text{Ce}$ (Ce-activated LaBr_3) detectors positioned in a planar fashion symmetrically around the ^{137}Cs γ -ray source (Fig. 1a). Thus, the angle between two adjacent detectors amounted to 72° , and two groups of five detector pairs can be formed with relative angles of 72° (72° -group) and 144° (144° -group). The undesired Compton scattering of 661.66-keV γ -rays between two detectors was suppressed by thick lead shields. Data were recorded for a measurement time of 52.7 days.

The spectrum of the time-difference (Δt) between two hits in two detectors of the 72° -group is presented in Fig. 1b inset. The prompt time-coincidence peak centred at $\Delta t = 0$ sits on a flat background caused by random coincidences. The time condition $|\Delta t| \leq 1.2$ ns selects true (prompt) as well as random coincidences, resulting in the energy-sum spectrum $E_1 + E_2$ shown in Fig. 1b (orange filled circles). The contribution of random coincidences is determined by imposing a wide time gate on the flat background and scaling the corresponding energy-sum spectrum with a factor of 0.0214 to correct for the different widths of both conditions (Fig. 1b, green filled triangles). The random coincidences dominate the statistical uncertainty.

The final energy-sum spectrum of the 72° -group—after subtraction of random coincidences—is shown in Fig. 2a. We imposed the energy condition $|E_1 - E_2| < 300$ keV to obtain further background suppression (see Supplementary Information). The spectrum is well described by assuming a superposition of an exponential background and a Gaussian peak at 661.6(1.6) keV in agreement with the expected energy sum (here and elsewhere, numerals in parentheses indicate the 1σ uncertainty). It has an area of 693(95) counts, establishing the observation of the competitive $\gamma\gamma$ -decay.

Reasonable explanations other than $\gamma\gamma/\gamma$ -decay can be excluded. First, let us assume that despite our careful assembly of the lead shields, the peak is caused by Compton-scattered 661.66-keV γ -rays (due, for example, to a small hole in the absorptive lead shields). The additional flight path of the γ quantum from one detector to the other would cause a time delay to which the $\text{LaBr}_3:\text{Ce}$ detectors are sensitive. Hence,

¹Institut für Kernphysik, Technische Universität Darmstadt, 64289 Darmstadt, Germany. ²Department of Physics, University of Saskatchewan, Saskatoon S7N5E2, Canada.

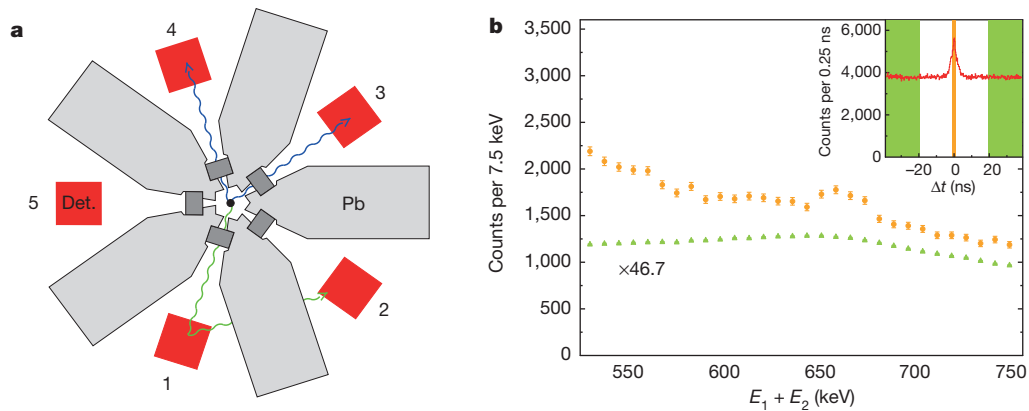


Figure 1 | The two main experimental obstacles to measuring the $\gamma\gamma/\gamma$ -decay. **a**, Schematic illustration of the experimental set-up. A 661.66-keV γ -ray (green wavy line) emitted from the central source (black dot) can deposit part of its energy in detector 1 (labelled red square) and then be scattered into detector 2 where it deposits its remaining energy. Lead (grey) is used to suppress this effect relative to the $\gamma\gamma$ -events (blue wavy lines). Lead collimators (dark grey) are used for background suppression. **b**, Energy-sum spectra of the 72° -group with time gates on the prompt coincidence peak (orange filled

circles; $|\Delta t| \leq 1.2$ ns) and on the random coincidences (green filled triangles; $20 \text{ ns} \leq |\Delta t| \leq 76$ ns). The ordinate values should be multiplied by the given factor. Inset, time difference spectrum (red trace) and parts of the imposed time gates (orange bar (prompt) and green areas (random)) used to obtain the energy-sum spectra of **b**. The energy condition $|E_1 - E_2| < 300$ keV was imposed on data reported in **b** as well as on that shown in the inset. Error bars, ± 1 s.d.

this possibility can be excluded by investigating the shape of the time-difference spectra. A narrow energy gate on the $\gamma\gamma/\gamma$ -sum-energy peak results in the time-difference spectrum presented in Fig. 2b (orange filled circles). The contribution of the background counts (solid green curve) is determined with energy gates on the background. The time-difference spectrum follows a Gaussian distribution centred at $\Delta t = 0$ (solid orange curve) with a full-width at half-maximum that is in agreement with the time resolution of our LaBr₃:Ce detectors, namely 1 ns. In the case of Compton scattering, a superposition of two Gaussians with centroids at $\Delta t = \pm 0.8$ ns (solid red curve in Fig. 2b) would be expected, in contradiction to the observed time spectrum.

A second possibility for the origin of the peak is a sequential decay via the $1/2^+$ state at an excitation energy of 283.5 keV (Fig. 2a inset), corresponding to two subsequent γ -decays. Our observed γ -ray energy spectra of individual photons—gated on $\gamma\gamma/\gamma$ -events—are continuous and do not peak at the transition energies to and from the $1/2^+$ state, thereby excluding sequential decay. In addition, the

γ -branch to the $1/2^+$ state has been measured recently and amounts to only $1.12(9) \times 10^{-7}$ (ref. 14), less than 6% of our measured $\gamma\gamma/\gamma$ -branching ratio.

In addition, data were recorded for the detector group forming opening angles of 144° . A Gaussian peak is found in the corresponding energy-sum spectrum at an energy of 664.2(28) keV and with an area of 307(78) counts (Fig. 3). This shows that the $\gamma\gamma/\gamma$ -decay of the $11/2^-$ isomer of ^{137}Ba exhibits a very pronounced angular correlation of the two emitted γ quanta, characteristic of the multiplicities involved in the transition.

In the theoretical treatment of the $\gamma\gamma$ -decay process, it follows from second-order perturbation theory that the calculation of the $\gamma\gamma$ -decay transition rate^{3,15} involves a summation over a complete set of states, subject to well known electromagnetic selection rules based on the multipolarity and character of the transition. In a reasonable approximation, the possible decay paths can be restricted to two cases (see Supplementary Information for a justification of this approximation), which we will describe here in the language of

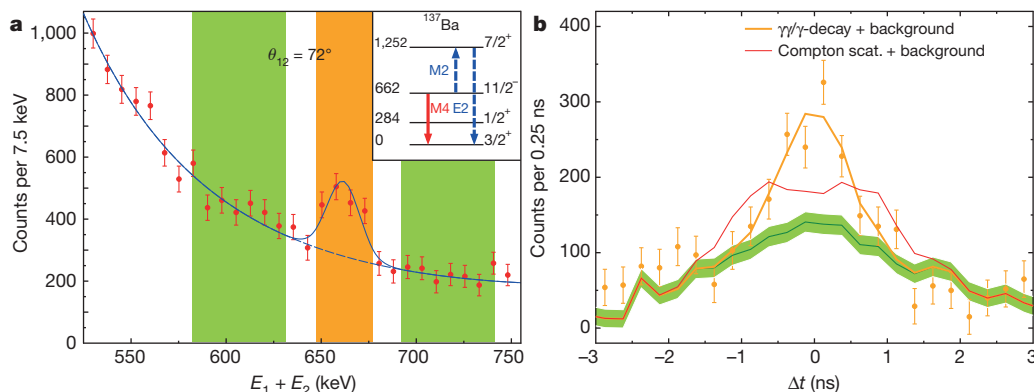


Figure 2 | Energy-sum spectrum and energy-gated time spectra of the 72° -group. **a**, Energy-sum spectrum $E_1 + E_2$ after subtraction of the random coincidences (requiring the energy condition $|E_1 - E_2| < 300$ keV). The spectrum is fitted with a superposition of a Gaussian and an exponential function to describe the peak and the background, respectively. The orange ($647 \text{ keV} < E_1 + E_2 < 677 \text{ keV}$) and green areas ($582 \text{ keV} < E_1 + E_2 < 632 \text{ keV}$ and $692 \text{ keV} < E_1 + E_2 < 742 \text{ keV}$) represent the energy conditions employed to obtain the time-difference spectra displayed in **b**. Inset, level scheme of ^{137}Ba with level energies in keV on the left and spin-parity J^π quantum

numbers on the right. The solid red arrow labels the single- γ transition of the $11/2^-$ state to the ground state. The dashed blue arrows illustrate the structure of the matrix element of the $\gamma\gamma/\gamma$ -decay involving the lowest $7/2^+$ state. **b**, Time difference spectra. The orange data points correspond to the orange region in **a**, while the green solid line with surrounding shading shows the background corresponding to the green area in **a**. The solid orange curve shows the expected time spectrum for $\gamma\gamma$ -decay, while the solid red curve shows the expected time spectrum, assuming the peak at 661.66 keV was caused by Compton-scattered γ -rays. Error bars in **a**, **b** are ± 1 s.d.

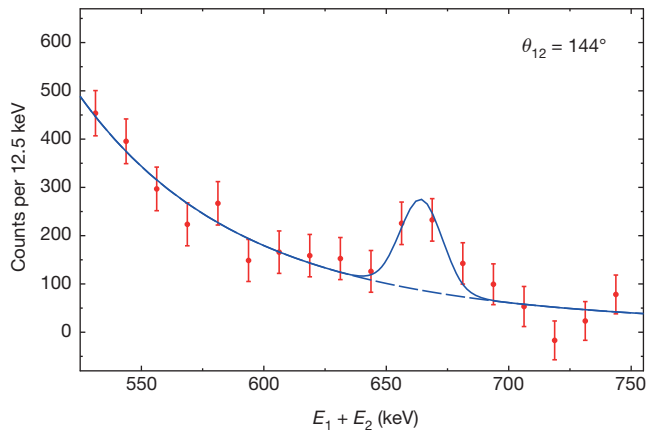


Figure 3 | Energy-sum spectrum of the 144°-group. The energy-sum spectrum $E_1 + E_2$ of the 144°-group after subtraction of the random coincidences (with the energy condition $|E_1 - E_2| < 250$ keV) is shown by the red data points. The Gaussian (solid blue curve) on top of an exponential background (dashed-blue curve) with a peak area of 307(78) counts is located at an energy of 664.2(28) keV, in agreement with the transition energy of 661.66 keV. Error bars, ± 1 s.d.

‘Feynman diagrams’: in the first case, the nucleus undergoes virtual M2 transitions from the $11/2^-$ state to intermediate $7/2^+$ states (Fig. 2a inset), which are coupled by E2 transitions to the $3/2^+$ ground state. The final matrix element is obtained by a summation over all intermediate $7/2^+$ states of ^{137}Ba , whose contributions add up coherently. The second possibility is an E3M1 transition through intermediate $5/2^+$ states. The theory of the $\gamma\gamma$ -decay is fully developed (see, for example, the Appendix of ref. 3; a first application of the formalism to ^{137}Ba was performed by D.J. Millener (personal communication)), and the differential branching ratio for the current case is approximately given by

$$\frac{d^5\Gamma_{\gamma\gamma}}{d\omega d\Omega d\Omega'} = A_{\text{qq}}(\alpha_{\text{E2M2}}^2, s) + A_{\text{od}}(\alpha_{\text{M1E3}}^2, s) + A_{\text{x}}(\alpha_{\text{E2M2}}\alpha_{\text{M1E3}}, s) \quad (1)$$

where $\Gamma_{\gamma\gamma}$ is the decay width for two-photon emission, and s stands for $\{L, L', I_n, \omega, \theta_{12}\}$; here L and L' label the multipoles of the two virtual transitions involved, I_n is the spin of the intermediate state n , ω is the energy of one of the γ -rays, θ_{12} denotes the angle between the two emitted γ -rays (that is, in our experiment $\theta_{12} = 72^\circ$ and 144°) and $d\Omega$ and $d\Omega'$ are their solid angle elements. The coefficients α_{E2M2} and α_{M1E3} contain the nuclear structure information and are defined as

$$\alpha_{S'L'SL} = \sum_n \frac{\langle \frac{3}{2}^+ \| S'L' \| I_n \rangle \cdot \langle I_n \| SL \| \frac{11}{2}^- \rangle}{E_n - 0.5E_0} \quad (2)$$

where S and S' are the multipole characters. Each intermediate state $|I_n\rangle$ contributes with a product of matrix elements connecting it to the ground state ($\langle 3/2^+ \| S'L' \| I_n \rangle$) and to the $11/2^-$ isomer ($\langle I_n \| SL \| 11/2^- \rangle$) weighted approximately with the inverse of its excitation energy E_n . The quadrupole–quadrupole (A_{qq}) and octupole–dipole (A_{od}) terms of equation (1) depend only on α_{E2M2}^2 and α_{M1E3}^2 , respectively, while the interference term A_{x} depends on their product. Each term has a characteristic dependence on θ_{12} and the γ -ray energy ω . Thus, the values of α_{E2M2} , α_{M1E3} and the sign of $\alpha_{\text{E2M2}}\alpha_{\text{M1E3}}$ can be determined from an analysis of the measured angular correlation and the energy spectra of the emitted γ -rays. More information and the full form of equation (1) is given in Supplementary Information.

In Fig. 4, the result of a simultaneous fit to the energy spectrum (Fig. 4a) and angular correlation (Fig. 4b) is presented. The available data are well described by equation (1) (solid blue curve), and the contribution of the A_{qq} term (dashed blue curve) dominates in com-

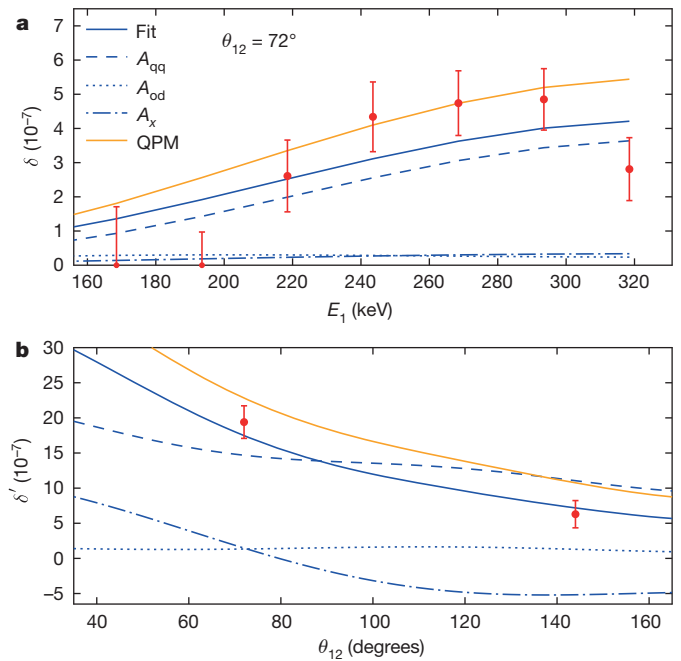


Figure 4 | Measured energy spectrum and angular correlation. δ (δ') is the differential $\gamma\gamma$ -decay width from equation (1) integrated over the corresponding energy bin, divided by the γ -decay width Γ_{γ} and evaluated at the average emission angles $\theta_{12} = 72^\circ$ and 144° (see Supplementary Information equation (2) for details of δ and δ'). **a**, Data points, dependence of δ on the energy E_1 of the lower-energy γ -ray for $\theta_{12} = 72^\circ$. Also shown are the results of the fit of equation (1) (solid blue curve) and the contributions of the A_{qq} (dashed blue curve), A_{od} (blue dotted curve) and A_{x} terms (blue dash-dotted curve). Results of the QPM calculation (orange solid curve) are compared to the data. **b**, Data points, dependence of δ' on θ_{12} integrated over the energy bin $|E_1 - E_2| < 250$ keV. Other curves as in **a**. All values are corrected for the small contribution of the cascade of the decay via the $1/2^+$ state. Error bars, ± 1 s.d.

parison to the A_{od} term (dotted blue curve). The angular correlation is very sensitive to the product of α_{E2M2} and α_{M1E3} . The final branching ratio amounts to $\Gamma_{\gamma\gamma}/\Gamma_{\gamma} = 2.05(37) \times 10^{-6}$, and the values for the two matrix elements are given in Table 1.

In Fig. 4 and Table 1 our data are also compared to predictions of the microscopic quasiparticle–phonon model (QPM)¹⁶ (solid orange curves). Further details of the QPM calculation are given in Supplementary Information. The QPM reproduces the experimental branching ratio and the two matrix elements rather well. As experimentally observed, it predicts a dominant A_{qq} term and a positive sign for $\alpha_{\text{E2M2}}\alpha_{\text{M1E3}}$. According to the QPM results, the value of α_{E2M2} is dominated by the matrix element coupled to the lowest $7/2^+$ state at an excitation energy of 1,252 keV (Supplementary Information). The reason is twofold: the lowest $7/2^+$ is close in energy to the $11/2^-$ state, reducing the size of the denominator in equation (2), and almost all of the low-lying E2-strength is concentrated in the $3/2^+ \rightarrow 7/2_1^+$ transition, resulting in a large $\langle 3/2^+ \| E2 \| 7/2_1^+ \rangle$ matrix element.

In conclusion, we have observed the $\gamma\gamma$ -decay of a nuclear transition in competition with an allowed γ -decay. Our results demonstrate that theory provides a realistic description of the double-photon decay, and that the present state of the art of experimental equipment allows such a process to be measured even in the presence of the direct one-photon

Table 1 | Measured and theoretical parameter values

Parameter	Experiment	QPM theory
$\Gamma_{\gamma\gamma}/\Gamma_{\gamma} (10^{-6})$	2.05(37)	2.69
$\alpha_{\text{E2M2}} (e^2 \text{ fm}^4 \text{ MeV}^{-1})$	+33.9(2.8)	+42.60
$\alpha_{\text{M1E3}} (e^2 \text{ fm}^4 \text{ MeV}^{-1})$	+10.1(4.2)	+9.50

See text for details of parameters. The uncertainties include the statistical error from the fit (± 1 s.d.) and systematic contributions.

decay. This opens a new field of studies that were not previously possible. While the branching ratio is probably too small for the observation of the $\gamma\gamma$ -decay of a first excited 2^+ state, which is present in nearly all even-even nuclei, we expect that there are other odd- A nuclei where the competitive $\gamma\gamma$ -decay can be measured. Such experiments will be difficult and time-consuming, and the number of suitable nuclides obtainable in an isomeric state is limited. Nevertheless, the sensitivity of the experiment is now sufficient to enable investigation of so-far-unexplored nuclear structure observables, namely the (off-diagonal) nuclear electric and magnetic polarizabilities.

Received 13 May; accepted 27 August 2015.

1. Watson, B. A., Bardin, T. T., Becker, J. A. & Fisher, T. R. Two-photon decay of the 6.05-MeV state of ^{16}O . *Phys. Rev. Lett.* **35**, 1333–1336 (1975).
2. Schirmer, J. *et al.* Double gamma decay in ^{40}Ca and ^{90}Zr . *Phys. Rev. Lett.* **53**, 1897–1900 (1984).
3. Kramp, J. *et al.* Nuclear two-photon decay in $0^+ \rightarrow 0^+$ transitions. *Nucl. Phys. A* **474**, 412–450 (1987).
4. Mokler, P. H. & Dunford, R. W. Two-photon decay in heavy atoms and ions. *Phys. Scr.* **69**, C1–C9 (2004).
5. Ilakovac, K., Uroic, M., Majer, M., Pasic, S. & Vukovic, B. Two-photon decay of k-shell vacancy states in heavy atoms. *Radiat. Phys. Chem.* **75**, 1451–1460 (2006).
6. Göppert, M. Über die Wahrscheinlichkeit des Zusammenwirkens zweier Lichtquanten in einem Elementarakt. *Naturwissenschaften* **17**, 932 (1929).
7. Göppert-Mayer, M. Über Elementarakte mit zwei Quantensprüngen. *Ann. Phys.* **401**, 273–294 (1931).
8. Göppert-Mayer, M. Double beta-disintegration. *Phys. Rev.* **48**, 512–516 (1935).
9. Hayes, A. C. *et al.* Two-photon decay of the first excited 0^+ state in ^{16}O . *Phys. Rev. C* **41**, 1727–1735 (1990).
10. Beusch, W. Über Zweiquanten-Übergänge an Ba^{137} . *Helv. Phys. Acta* **33**, 363–394 (1960).
11. Basenko, V. K., Berlizov, A. N. & Prokopets, G. A. Estimation of the probability of two-photon decay of the 0.662 MeV ^{137}Ba state. *Bull. Russ. Acad. Sci. Phys.* **56**, 94 (1992).
12. Millener, D. J., Sutter, R. J. & Alburger, D. E. 2-gamma decay of the 662-keV isomer in ^{137}Ba . *Bull. Am. Phys. Soc.* **56**, DNP.CF.8 (2011); available at <http://meetings.aps.org/link/BAPS.2011.DNP.CF.8> (2011).
13. Lister, C. J. *et al.* A search for 2-photon emission from the 662 keV state in ^{137}Ba . *Bull. Am. Phys. Soc.* **58**, DNP.CE.3 (2013); available at <http://meetings.aps.org/link/BAPS.2013.DNP.CE.3> (2013).
14. Moran, K. *et al.* E5 decay from the $J^\pi = 11/2^-$ isomer in ^{137}Ba . *Phys. Rev. C* **90**, 041303 (2014).
15. Friar, J. L. Low-energy theorems for nuclear Compton and Raman scattering and $0^+ \rightarrow 0^+$ two-photon decays in nuclei. *Ann. Phys.* **95**, 170–201 (1975).
16. Soloviev, V. G. *Theory of Atomic Nuclei: Quasiparticles and Phonons* (Institute of Physics Publishing, Bristol, 1992).

Supplementary Information is available in the online version of the paper.

Acknowledgements This work was supported by the State of Hesse under the Helmholtz International Center for FAIR (HIC for FAIR) and by the German Research Council (DFG) under grant no. SFB 634. We thank D. J. Millener, R. J. Sutter for an open discussion and for generously sharing their preliminary results before publication, and C. J. Lister for discussions. H.S. thanks Dirk Schwalm for raising interest in this topic.

Author Contributions C.W. and H.S. performed the data analysis and derived the equations given in Supplementary Information. C.W., H.S. and N.P. contributed to the interpretation of the results. C.W. and R.L. were responsible for the set-up of the $\text{LaBr}_3:\text{Ce}$ detector array and the data acquisition system. V.Yu.P. performed the QPM calculation. C.W., H.S., N.P. and T.A. prepared the manuscript. All authors discussed the results, commented on and contributed to the manuscript.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to H.S. (hscheit@ikp.tu-darmstadt.de).

A two-qubit logic gate in silicon

M. Veldhorst¹, C. H. Yang¹, J. C. C. Hwang¹, W. Huang¹, J. P. Dehollain¹, J. T. Muhonen¹, S. Simmons¹, A. Laucht¹, F. E. Hudson¹, K. M. Itoh², A. Morello¹ & A. S. Dzurak¹

Quantum computation requires qubits that can be coupled in a scalable manner, together with universal and high-fidelity one- and two-qubit logic gates^{1,2}. Many physical realizations of qubits exist, including single photons³, trapped ions⁴, superconducting circuits⁵, single defects or atoms in diamond^{6,7} and silicon⁸, and semiconductor quantum dots⁹, with single-qubit fidelities that exceed the stringent thresholds required for fault-tolerant quantum computing¹⁰. Despite this, high-fidelity two-qubit gates in the solid state that can be manufactured using standard lithographic techniques have so far been limited to superconducting qubits⁵, owing to the difficulties of coupling qubits and dephasing in semiconductor systems^{11–13}. Here we present a two-qubit logic gate, which uses single spins in isotopically enriched silicon¹⁴ and is realized by performing single- and two-qubit operations in a quantum dot system using the exchange interaction, as envisaged in the Loss–DiVincenzo proposal². We realize CNOT gates via controlled-phase operations combined with single-qubit operations. Direct gate-voltage control provides single-qubit addressability, together with a switchable exchange interaction that is used in the two-qubit controlled-phase gate. By independently reading out both qubits, we measure clear anticorrelations in the two-spin probabilities of the CNOT gate.

Quantum dots have long been considered an attractive physical platform for quantum information processing². Large arrays can be conveniently realized using conventional lithographic approaches. Initialization, read-out, control and coupling can be achieved through local electrical pulses, possibly in combination with magnetic resonance techniques. Early research focused mainly on III–V semiconductor compounds such as GaAs, resulting in single-spin qubits¹⁵, singlet-triplet qubits¹⁶ and exchange-only qubits¹⁷, which can be coupled capacitively¹¹ or via the exchange interaction^{12,13}. Although these approaches demonstrate the potential of quantum-dot qubits, strong dephasing due to the nuclear spin background has limited the quality of the quantum operations. A marked improvement in coherence times has been observed by defining the quantum dots in silicon^{9,18,19}, which can be isotopically purified¹⁴, such that quantum dots with single-spin fidelities above the threshold of surface codes¹⁰ can be realized⁹.

A scalable approach towards quantum computation ideally requires that the coupling between qubits can be turned on and off², so that single- and two-qubit operations can be selectively chosen. Here, we push silicon-based quantum information processing beyond the single-qubit level by realizing a controlled-phase (CZ) gate, which is commonly used in superconducting qubits⁵ and has been theoretically discussed for quantum-dot systems²⁰. This two-qubit gate, together with single-qubit gates, provides all of the necessary operations for universal quantum computation. In our system, each qubit is defined by the spin state of a single electron, with energies split by a large magnetic field of strength $B_0 = 1.4$ T. The single-qubit states are manipulated using spin-resonance techniques, through the local application of an oscillating magnetic field produced by an on-chip electron spin resonance (ESR) line. By exploiting the Stark shift, we electrically control the effective

g -factor of the qubits, to tune the Zeeman energy $E_Z = g\mu_B B_0$ and the associated qubit resonance frequency $\nu = E_Z/h$ for selective qubit control⁹, where μ_B is the Bohr magneton and h is the Planck constant. The two-qubit gate is then realized using electrical pulses that control the exchange coupling between the qubits.

Figure 1a shows a schematic and Fig. 1b shows a scanning electron microscope image of the double-quantum-dot structure fabricated on a ²⁸Si epilayer with a residual ²⁹Si concentration of 800 p.p.m. (ref. 14). The device consists of three aluminium layers, nine aluminium gates, an aluminium lead for ESR control²¹, and source, drain and reservoir leads that connect to a gate-induced two-dimensional electron gas using multilevel gate-stack silicon metal–oxide–semiconductor technology²². A single-electron transistor (SET) is formed to monitor the charge state of the quantum-dot system and for spin read-out using spin-to-charge conversion²³. For both our single- and two-qubit experiments, we tune the gate voltages to the appropriate dc operating regime and then adjust gate G_1 for qubit read-out, initialization and control.

Figure 1c shows the stability diagram of the double-quantum-dot system with charge occupancy (N_2, N_1) . The charge transitions of quantum dots D_1 and D_2 , which are underneath gates G_1 and G_2 , respectively, are distinguished by their gate voltage dependence and their capacitive coupling to the SET. We define qubit Q_1 by loading a single electron into D_1 , so that $N_1 = 1$; similarly for qubit Q_2 , we have $N_2 = 1$.

To characterize the individual qubits, we bias the gate voltages such that the tunnel time of the respective qubit to the reservoir is approximately 100 μ s and both qubits are measured in the $(1, 1)$ charge state. Clear Rabi oscillations are observed as a function of microwave pulse time τ_p for both qubits, as shown in Fig. 1d. Q_1 has a dephasing time of $T_2^* = 120$ μ s, with a coherence time T_2 that can be extended up to 28 ms using CPMG (Carr–Purcell–Meiboom–Gill) pulses⁹, and Q_2 has a dephasing time of $T_2^* = 61$ μ s (see Supplementary Information section 3).

We couple the qubits via the exchange interaction, as discussed in ref. 2, with an exchange coupling that is electrically controlled via the detuning energy ϵ (see Fig. 2a, b). We control the system in the $(1, 1)$ region, and read out Q_2 at the $(1, 1)$ – $(0, 1)$ transition and Q_1 at the $(0, 1)$ – $(0, 0)$ transition.

The presence of a sharp interface and large perpendicular electric fields increases the energy of all excited states²⁴, and allows us to consider only the lowest five energy states¹⁶ (see also Supplementary Information section 2). We can consequently describe the system in the rotating-wave approximation in the basis $[|Q_2, Q_1\rangle, \Psi = [|\uparrow, \uparrow\rangle, |\uparrow, \downarrow\rangle, |\downarrow, \uparrow\rangle, |\downarrow, \downarrow\rangle, |0, 2\rangle]$, with the effective Hamiltonian

$$H = \begin{bmatrix} \overline{E_Z} - \nu & \Omega & \Omega & 0 & 0 \\ \Omega & \delta E_Z/2 & 0 & \Omega & t_0 \\ \Omega & 0 & -\delta E_Z/2 & \Omega & -t_0 \\ 0 & \Omega & \Omega & -\overline{E_Z} + \nu & 0 \\ 0 & t_0 & -t_0 & 0 & U - \epsilon \end{bmatrix} \quad (1)$$

¹Centre for Quantum Computation and Communication Technology, School of Electrical Engineering and Telecommunications, The University of New South Wales, Sydney, New South Wales 2052, Australia. ²School of Fundamental Science and Technology, Keio University, 3-14-1 Hiyoshi, Kohoku-ku, Yokohama 223-8522, Japan.

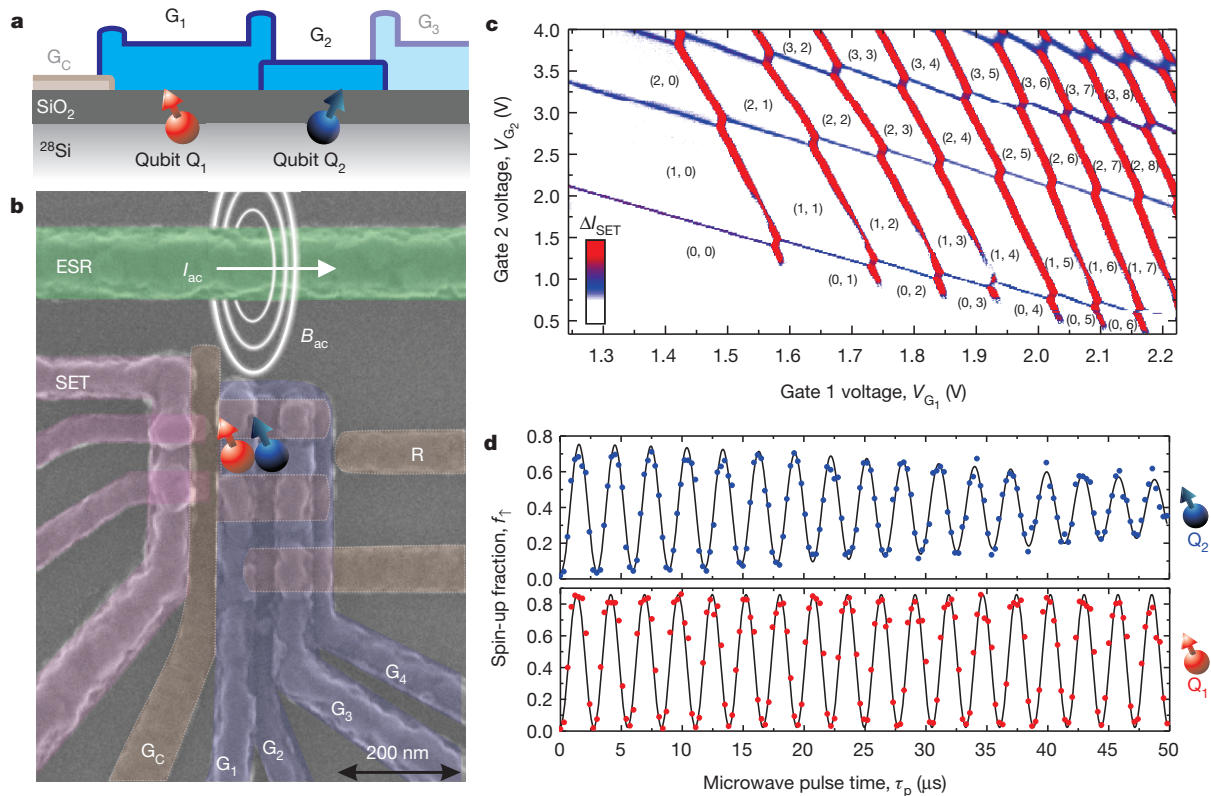


Figure 1 | Silicon two-qubit logic device, incorporating SET read-out and selective qubit control. **a**, **b**, Schematic (**a**) and scanning electron microscope coloured image (**b**) of the device. The quantum dot structure (labels G_C and G_{1-4}) can be operated as a single or double quantum dot by appropriate biasing of gate electrodes G_1 – G_4 , where we choose here to confine the dots $D_{1,2}$ underneath gates $G_{1,2}$, respectively. The confinement gate G_C runs underneath G_1 – G_3 and confines the quantum dot on all sides except on the reservoir (R) side. Qubit operation is achieved via an ac current I_{ac} through the ESR line, resulting in an ac magnetic field B_{ac} . **c**, Stability diagram of the double quantum dot obtained by monitoring the current I_{SET} through the capacitively coupled SET. The numbers in parentheses are the change occupancies of $D_{2,1}$: (N_2, N_1) .

The difference in distance to the SET results in different capacitive coupling, such that the individual dots can be easily distinguished. The tunnel coupling of the fourth transition ($N_1 = 3 \rightarrow 4$) of D_1 is relatively weak, which is due to valley and spin filling, because there is only one state in the lowest orbital that can be occupied. Q_1 and Q_2 are realized by depleting D_1 and D_2 to the last electron. **d**, The quantum dot qubits can be individually controlled by electrically tuning the ESR resonance frequency using the Stark shift⁹. Clear Rabi oscillations for both qubits are observed. All measurements were performed in a dilution refrigerator with base temperature $T \approx 50$ mK and a dc magnetic field of strength $B_0 = 1.4$ T.

where $\overline{E_Z}$ is the mean Zeeman energy, δE_Z is the difference in Zeeman energy between the dots, Ω is the Rabi frequency, ν is the microwave frequency and t_0 is the tunnel coupling; for simplicity we have scaled the system such that $\hbar = 1$. In the experiments, we control ϵ by fast pulsing only on G_1 . The read-out on Q_2 (R_2) and control (C) bias points are depicted in Fig. 2a; we pulse close to the $(1, 1) \rightarrow (1, 2)$ transition, which has the same energy level structure as the $(1, 1) \rightarrow (0, 2)$ transition, shown in Fig. 2b. Single-qubit operations are realized with Rabi frequency Ω , by matching ν to the resonance frequency of one of the qubits. The presence of exchange coupling between the qubits alters the Zeeman levels as shown in Fig. 2b, where the finite coupling between the qubits causes an anticrossing between the $|0, 2\rangle$ and $|1, 1\rangle$ states. We experimentally map out the energy levels in the vicinity of the anticrossing, as shown in Fig. 2c, d.

In Fig. 2c, we have initialized Q_1 and Q_2 to spin down and, by applying a π -pulse to Q_2 (π_{X,Q_2}), we map out the resonance frequency of Q_2 as a function of detuning. We measure exchange couplings of more than 10 MHz, above which T_2^* of Q_2 becomes shorter than the π -pulse time $\tau_\pi = 1.5$ μ s and so spin flips cannot occur.

Initialization of antiparallel spin states is possible by pulsing to the $(1, 2)$ and returning to the $(1, 1)$ charge states (labelled I_{AP} in Fig. 2a). In this sequence, an electron tunnels from D_2 to D_1 (I in Fig. 2d) followed by an electron tunnelling from the reservoir R to D_2 (II in Fig. 2d). After returning to $(1, 1)$, the electron from D_2 tunnels back to R (III in Fig. 2d), and one of the two electrons on D_1 , which are in a

singlet state, tunnels to D_2 (IV in Fig. 2d). With this initialization into an antiparallel spin state, when we apply a microwave pulse on Q_2 , the spin-up fraction f_1 approaches 1/2, except when the microwave frequency matches a resonance frequency of Q_2 . Owing to the finite exchange interaction, there are two resonance frequencies. The lower frequency rotates the antiparallel state towards a combination of $|\downarrow, \downarrow\rangle$ and $|\downarrow, \uparrow\rangle$, where Q_2 always ends up as spin down. The higher frequency rotates the antiparallel state towards a combination of $|\uparrow, \downarrow\rangle$ and $|\uparrow, \uparrow\rangle$, where Q_2 always ends up as spin up. The results are depicted in Fig. 2d, which shows a decrease in f_1 at the lower branch and an increase of f_1 at the upper branch, demonstrating an exchange spin funnel, where both branches are visible (see Supplementary Information section 4 for further details). The two-qubit gate is conveniently realized using a quantum dot CZ gate²⁰ (see Supplementary Information section 6 for theoretical details). This approach allows individual control over the qubits in the absence of interaction and its associated noise, while using the coupling to perform two-qubit operations with a frequency that can be much higher than the single-qubit Rabi rotation frequency.

As described by equation (1) and depicted in Fig. 2b, changing ϵ modifies the qubit resonance frequencies of Q_1 and Q_2 and introduces an effective detuning frequency $\nu_{\uparrow\downarrow,(\uparrow\downarrow)}$, such that one qubit acquires a time-integrated phase shift $\phi_{\uparrow\downarrow,(\uparrow\downarrow)}$ that depends on the \hat{z} component of the spin state of the other qubit, and vice versa. The exchange coupling and $\nu_{\uparrow\downarrow,(\uparrow\downarrow)}$ are maximized at the anticrossing $\epsilon = U$; see Fig. 2b. When

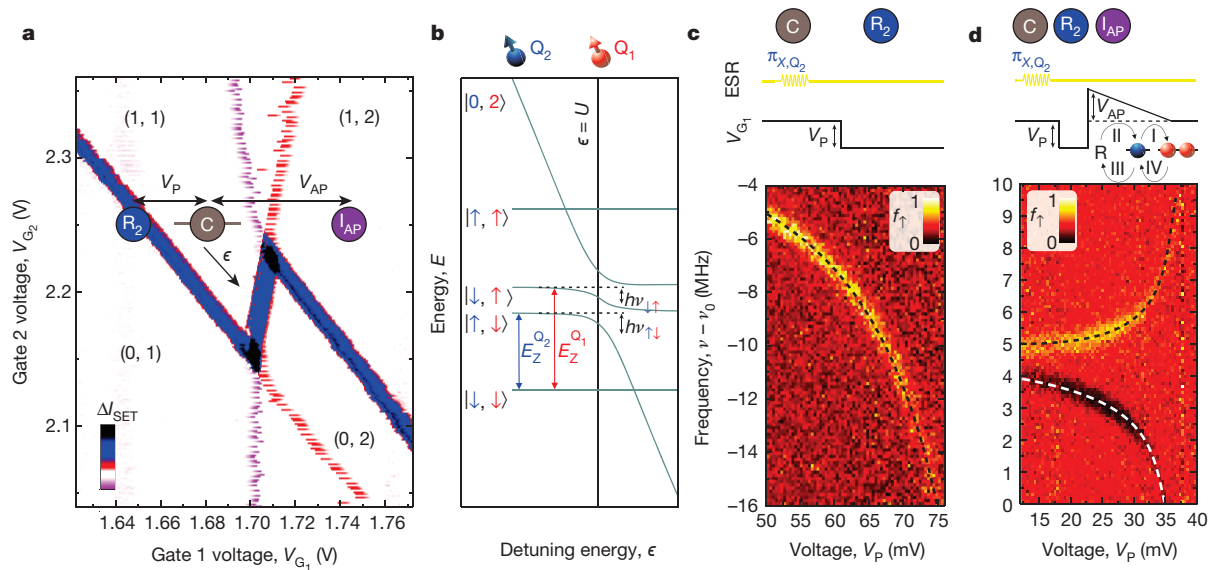


Figure 2 | Exchange spin funnel. **a**, Close up of the operation regime of the $(1, 1) \rightarrow (0, 2)$ charge states. We lowered the R–D₂ coupling so that the tunnelling time is approximately 100 μ s, matching the qubit experiments. In this range of weak R–D₁ coupling, the emptying and filling of D₁ is hysteretic with gate voltage, because the mutual charging energy becomes relevant²⁷, as D₁ can only tunnel when it aligns in energy with D₂. R₂ represents the read-out on Q₂; I_{AP} represents the antiparallel initialization. **b**, Schematic of the coupling between Q₁ and Q₂ using the exchange interaction at the $|1, 1\rangle \rightarrow |0, 2\rangle$ transition. By electrically tuning the g factors⁹ of Q₁ and Q₂, we control the individual qubit resonance frequencies over 10 MHz. Here, we tune to a frequency difference of

40 MHz (the difference is exaggerated in the schematic for clarity) for individual qubit control. **c**, ESR spectrum of the $|\downarrow, \downarrow\rangle \rightarrow |\uparrow, \downarrow\rangle$ transition as a function of increasing detuning. The data have been offset by a frequency $\nu_0 = 39.14$ GHz and the spin-up fractions are normalized for clarity. The dashed lines are fits using equation (1) and assuming $t_0 = 900$ MHz, a Stark shift of 19 MHz V^{-1} , and that the top gates have a lever arm of 0.2 eV V^{-1} . **d**, As for **c**, but with an additional pulse of amplitude V_{AP} (see schematic) so that we initialize antiparallel spin states (AP) and observe both the $|\uparrow, \downarrow\rangle \rightarrow |\downarrow, \downarrow\rangle$ and $|\downarrow, \uparrow\rangle \rightarrow |\uparrow, \uparrow\rangle$ transitions. The labels I–IV indicate electron tunnelling between R and D_{1,2}; see text for details.

a CZ operation is performed such that $\phi_{\uparrow\downarrow} + \phi_{\downarrow\uparrow} = \pi$, the operation differs only by an overall phase from the basis CZ gate²⁵. This overall phase can be removed using single-qubit pulses or via voltage pulses exploiting the Stark shift⁹. To realize a CNOT operation using the CZ gate, a CZ(π) rotation is performed in between two $\pi/2$ -pulses on Q₂ that have a phase difference $\phi_{\uparrow\downarrow}$.

Figure 3a shows the spin-up fraction f_{\uparrow} of Q₂ after applying a $(\pi/2)_X$ -pulse and a $(\pi/2)_Y$ -pulse on Q₂ separated by an interaction time τ_Z with increasing exchange coupling, set via ϵ and tuned by the voltage V_{CZ} . Plotting the frequency $\nu_{\uparrow\downarrow}$ as a function of V_{CZ} (Fig. 3b) gives a trend consistent with that observed via ESR mapping, as shown in Fig. 2c, d. The two-qubit dephasing time $T_{2,CZ}^*$ is the free induction decay time of

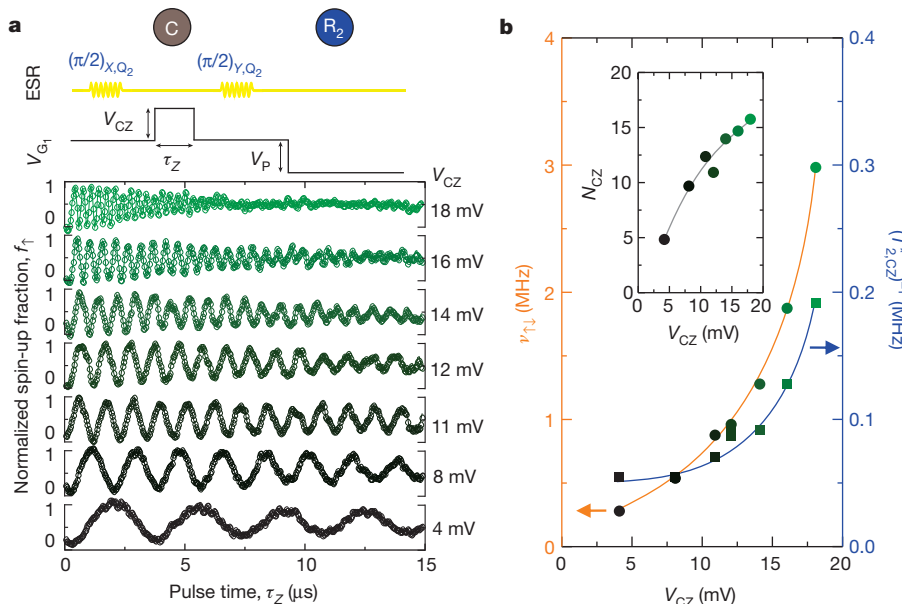


Figure 3 | Controlled phase (CZ) gate operation time. **a**, Spin-up fraction after applying a $(\pi/2)_X$ -pulse, a CZ operation and a $(\pi/2)_Y$ -pulse, for increasing qubit interaction. The exchange coupling is controlled via ϵ , set with V_{G1} , resulting in a tunable two-qubit operation frequency $\nu_{\uparrow\downarrow}$. **b**, By fitting the data in **a**, we map out $\nu_{\uparrow\downarrow}$ and $T_{2,CZ}^*$ as a function of V_{CZ} . The orange and blue colouring

and the arrows indicate the axis each data set corresponds to; the colouring of the data corresponds to that in **c**, indicating V_{CZ} . The inset shows the number of possible CZ rotations N_{CZ} . Although $T_{2,CZ}^*$ decreases with coupling, the number of possible two-qubit rotations continues to increase.

the two-qubit system. At large values of detuning ($\epsilon \rightarrow \infty$) or when the interaction vanishes ($t_0 \rightarrow 0$), $T_{2,CZ}^*$ reduces to the single-qubit Ramsey T_2^* . We obtain $T_{2,CZ}^*$ by fitting an exponential to the decay of the oscillations in Fig. 3a. These values of $T_{2,CZ}^*$ are plotted along with the measured $v_{\uparrow\downarrow}$ values in Fig. 3b. We find that the two-qubit dephasing rate $(T_{2,CZ}^*)^{-1}$ rises in step with the exchange coupling and $v_{\uparrow\downarrow}$, which is to be expected because $\delta v_{\uparrow\downarrow}/\delta V$ also increases with $v_{\uparrow\downarrow}$, meaning that the qubit system becomes increasingly sensitive to electrical noise²⁶. Despite this, the total number of oscillations $N_{CZ} = v_{\uparrow\downarrow} T_{2,CZ}^*$ also increases with $v_{\uparrow\downarrow}$, as shown in Fig. 3b. In Supplementary Information section 7, we show an optimized sequence where $T_{2,CZ}^* = 8.3 \mu\text{s}$ and $v_{\uparrow\downarrow} = 3.14 \text{ MHz}$, such that $N_{CZ} > 26$.

For all of the experiments described in Figs 2 and 3, we performed read-out only on Q_2 , owing to its proximity to the reservoir used for spin selective read-out. However, to demonstrate a two-qubit CNOT gate it is desirable to read out the state of both qubits, allowing the observation of non-classical correlations. This requires a more complex pulsing protocol on G_1 (Fig. 4a), to first read Q_2 at the charge transition $(0,0)-(0,1)$, subsequently read and initialize Q_1 at

the $(0,1)-(1,1)$ transition, and initialize Q_2 at the $(0,0)-(0,1)$ transition. The read-out ($R_{1,2}$) and control (C) bias points are shown on the charge stability map in Fig. 4b. Following two-qubit read-out of the previous state, the system is prepared as $|\downarrow, \downarrow\rangle$, after which single-qubit rotations are applied to each qubit to prepare any desired initial two-qubit state.

Figure 4c shows the measured states of Q_1 and Q_2 after applying the CZ gate as a function of τ_Z , with Q_1 initialized to $|\uparrow\rangle$ (top panel) and $|\downarrow\rangle$ (bottom panel). As expected, the control qubit Q_1 (red) is not perturbed when exchange is turned on because it is in a basis state; however, the target qubit Q_2 (blue) is initialized to $\frac{1}{\sqrt{2}}(|\uparrow\rangle + |\downarrow\rangle)$ with a $(\pi/2)_X$ -pulse on Q_2 and so rotates about the equator of the Bloch sphere when exchange is turned on for time τ_Z . The strength of the exchange coupling is set so that Q_2 rotates about the Bloch sphere at double the frequency for $Q_1 = |\downarrow\rangle$ than for $Q_1 = |\uparrow\rangle$; this is reflected in the final state of Q_2 plotted in Fig. 4c (see Supplementary Information section 8 for further details). A CZ gate is realized at $\tau_Z = 480 \text{ ns}$, when $\phi_{\uparrow\downarrow} + \phi_{\downarrow\uparrow} = \pi$, and this is converted to a CNOT gate (the target qubit is

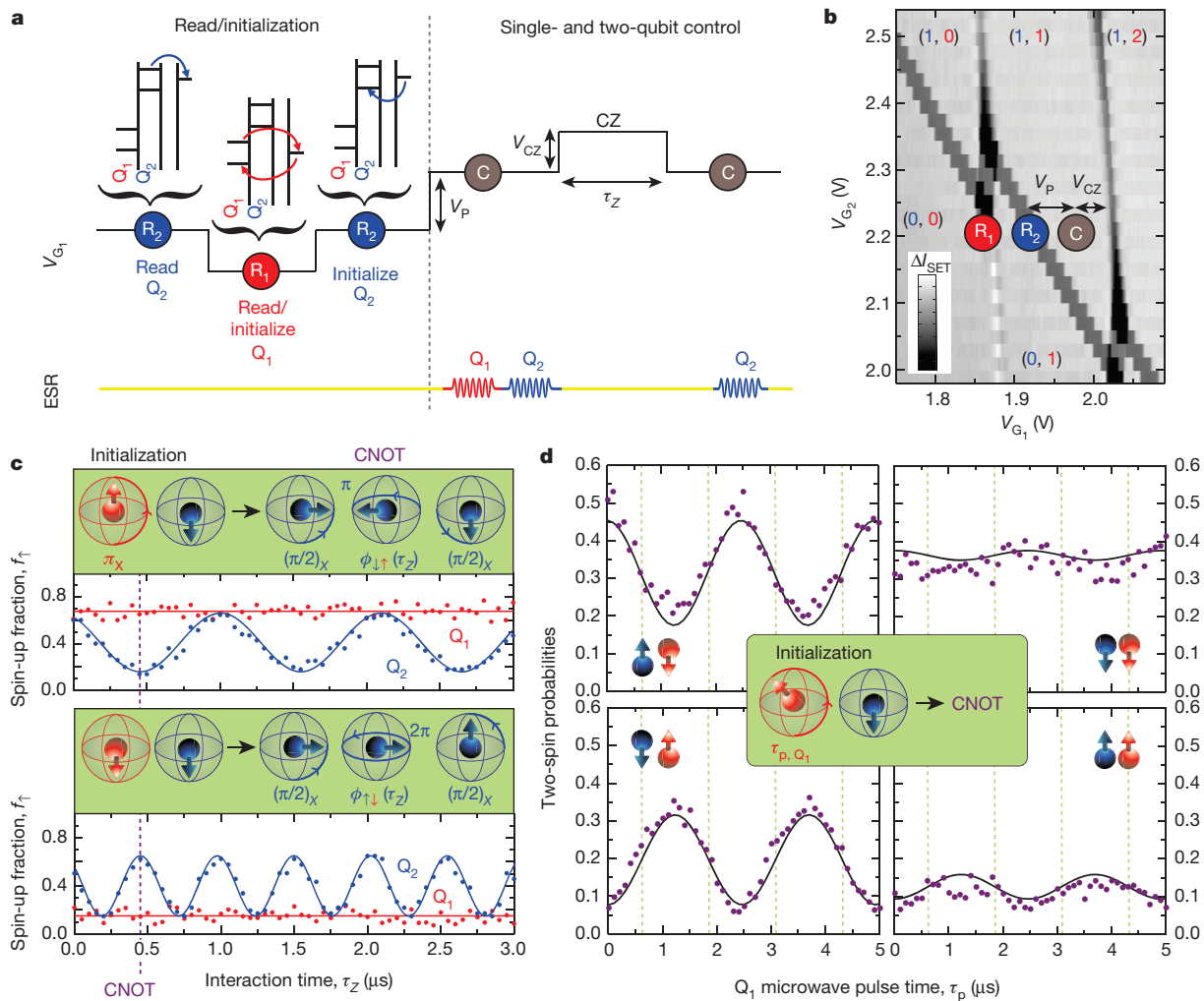


Figure 4 | Two-spin correlations for a two-qubit logic gate. **a**, Pulsing protocol for two-qubit read-out and single- and two-qubit operations. After read-out of Q_2 (R_2) and Q_1 (R_1), we pulse back to (R_2) to ensure proper initialization. Individual qubit operations are performed with high ϵ , whereas the CZ operation occurs in the presence of interaction. **b**, Stability diagram showing the operation regime. **c**, Spin-up fraction of both qubits after initializing Q_1 spin up (top) and spin down (bottom) using a microwave pulse and applying a controlled rotation using Q_2 as the target qubit. A CNOT gate is achieved in 480 ns, as indicated by the dotted purple line (see inset for the

corresponding Bloch sphere animation). **d**, Two-spin probabilities as functions of the microwave pulse length on Q_1 after applying a CNOT gate (see inset for the corresponding Bloch sphere animation), showing clear anticorrelations between the two qubit spin states. The different plots correspond to different spin states of $Q_{1,2}$, as indicated. The black lines correspond to fits based on a CNOT gate, and include the experimental read-out errors (see Supplementary Information section 9). The green dotted lines correspond to the intended maximally entangled states.

flipped when the control qubit is $|\downarrow\rangle$ by applying $(\pi/2)_X$ -pulses on Q_2 before and after the CZ.

We use the CNOT gate to create an entangled state of Q_1 and Q_2 . To realize this, we initialize the qubits first to the $|\downarrow, \downarrow\rangle$ state, then apply a varying microwave pulse time to rotate Q_1 into superposition states, with a Rabi time $\tau_{\text{Rabi}} = 2.4 \mu\text{s}$, and finally apply the CNOT gate. To demonstrate the CNOT gate, we convert the individual qubit spin-up fractions into two-spin probabilities; Fig. 4d shows the four possible two-spin probabilities. Clear oscillations are observed in the probabilities of the antiparallel states, $P(|\downarrow, \uparrow\rangle)$ and $P(|\uparrow, \downarrow\rangle)$, whereas these oscillations are almost absent in the probabilities of the parallel states, $P(|\downarrow, \downarrow\rangle)$ and $P(|\uparrow, \uparrow\rangle)$, thereby demonstrating the anticorrelations expected for the CNOT gate. The hints of oscillations in the symmetric spin states are probably due to read errors (which are included in the fitted line in Fig. 4d, see also Supplementary Information section 9); our current visibilities are not sufficient to demonstrate violation of the Bell inequality.

Future experiments will include improvements to the read-out fidelities, thus facilitating full two-qubit tomography. The qubit control fidelities could be further improved by lowering the sensitivity to electrical noise. Although these silicon qubits represent the smallest scalable two-qubit system reported so far, the complete fabrication process is compatible with standard CMOS (complementary metal-oxide-semiconductor) technology, and is also consistent with current transistor feature sizes, offering the prospect of realizing a large-scale quantum processor using the same silicon manufacturing technologies that have enabled the current information age.

Received 7 November 2014; accepted 22 July 2015.

Published online 5 October 2015.

- DiVincenzo, D. P. The physical implementation of quantum computation. *Fortschr. Phys.* **48**, 771–783 (2000).
- Loss, D. & DiVincenzo, D. P. Quantum computation with quantum dots. *Phys. Rev. A* **57**, 120–126 (1998).
- Kok, P. *et al.* Linear optical quantum computing with photonic qubits. *Rev. Mod. Phys.* **79**, 135–174 (2007).
- Brown, K. R. *et al.* Single-qubit-gate error below 10^{-4} in a trapped ion. *Phys. Rev. A* **84**, 030303 (2011).
- Barends, R. *et al.* Superconducting quantum circuits at the surface code threshold for fault tolerance. *Nature* **508**, 500–503 (2014).
- Waldherr, G. *et al.* Quantum error correction in a solid-state hybrid spin register. *Nature* **506**, 204–207 (2014).
- Dolde, F. *et al.* High-fidelity spin entanglement using optimal control. *Nature Commun.* **5**, 3371 (2014).
- Muhonen, J. T. *et al.* Storing quantum information for 30 seconds in a nanoelectronic device. *Nature Nanotechnol.* **9**, 986–991 (2014).
- Veldhorst, M. *et al.* An addressable quantum dot qubit with fault-tolerant fidelity. *Nature Nanotechnol.* **9**, 981–985 (2014).
- Fowler, A., Mariantoni, M., Martinis, J. M. & Cleland, A. N. Surface codes: towards practical large-scale quantum computation. *Phys. Rev. A* **86**, 032324 (2012).
- Shulman, M. D. *et al.* Demonstration of entanglement of electrostatically coupled singlet-triplet qubits. *Science* **336**, 202–205 (2012).
- Nowack, K. C. *et al.* Single-shot correlations and two-qubit gate of solid-state spins. *Science* **333**, 1269–1272 (2011).
- Brunner, R. *et al.* Two-qubit gate of combined single-spin rotation and interdot exchange in a double quantum dot. *Phys. Rev. Lett.* **107**, 146801 (2011).
- Itoh, K. M. & Watanabe, H. Isotope engineering of silicon and diamond for quantum computing and sensing applications. *Mater. Res. Soc. Commun.* **4**, 143–157 (2014).
- Koppens, F. H. L. *et al.* Driven coherent oscillations of a single electron spin in a quantum dot. *Nature* **442**, 766–771 (2006).
- Petta, J. R. *et al.* Coherent manipulation of coupled electron spins in semiconductor quantum dots. *Science* **309**, 2180–2184 (2005).
- Medford, J. *et al.* Self-consistent measurement and state tomography of an exchange-only spin qubit. *Nature Nanotechnol.* **8**, 654–659 (2013).
- Maune, B. M. *et al.* Coherent singlet-triplet oscillations in a silicon-based double quantum dot. *Nature* **481**, 344–347 (2012).
- Kawakami, E. *et al.* Electrical control of a long-lived spin qubit in a Si/SiGe quantum dot. *Nature Nanotechnol.* **9**, 666–670 (2014).
- Meunier, T., Calado, V. E. & Vandersypen, L. M. K. Efficient controlled-phase gate for single-spin qubits in quantum dots. *Phys. Rev. B* **83**, 121403(R) (2011).
- Dehollain, J. P. *et al.* Nanoscale broadband transmission lines for spin qubit control. *Nanotechnology* **24**, 015202 (2013).
- Angus, S. J., Ferguson, A. J., Dzurak, A. S. & Clark, R. G. Gate-defined quantum dots in intrinsic silicon. *Nano Lett.* **7**, 2051–2055 (2007).
- Elzerman, J. M. *et al.* Single-shot read-out of an individual electron spin in a quantum dot. *Nature* **430**, 431–435 (2004).
- Yang, C. H. *et al.* Spin-valley lifetimes in a silicon quantum dot with tunable valley splitting. *Nature Commun.* **4**, 2069 (2013).
- Ghosh, J. *et al.* High-fidelity controlled- σ^z gate for resonator-based superconducting quantum computers. *Phys. Rev. A* **87**, 022309 (2013).
- Dial, O. E. *et al.* Charge noise spectroscopy using coherent exchange oscillations in a singlet-triplet qubit. *Phys. Rev. Lett.* **110**, 146804 (2013).
- Yang, C. H. *et al.* Charge state hysteresis in semiconductor quantum dots. *Appl. Phys. Lett.* **105**, 183505 (2014).

Supplementary Information is available in the online version of the paper.

Acknowledgements We thank S. Bartlett for discussions and C. M. Cheng for contributions to the preparation of the experimental setup. We acknowledge support from the Australian Research Council (CE11E0001017), the US Army Research Office (W911NF-13-1-0024) and the NSW Node of the Australian National Fabrication Facility. M.V. acknowledges support from the Netherlands Organization for Scientific Research (NWO) through a Rubicon Grant. The work at Keio was supported in part by the Grant-in-Aid for Scientific Research by MEXT, in part by NanoQuine, in part by FIRST and in part by the JSPS Core-to-Core Program.

Author Contributions M.V., C.H.Y. and J.C.C.H. performed the experiments. M.V. and F.E.H. fabricated the devices. K.M.I. prepared and supplied the ^{28}Si epilayer wafer. W.H., J.P.D., J.T.M., S.S. and A.L. contributed to the preparation of the experiments. M.V., C.H.Y., A.M. and A.S.D. designed the experiment and discussed the results. M.V. analysed the results. M.V. and A.S.D. wrote the manuscript with input from all co-authors.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to M.V. (M.Veldhorst@unsw.edu.au) or A.S.D. (A.Dzurak@unsw.edu.au).

Peptoid nanosheets exhibit a new secondary-structure motif

Ranjan V. Mannige¹, Thomas K. Haxton¹, Caroline Proulx¹, Ellen J. Robertson¹, Alessia Battigelli¹, Glenn L. Butterfoss², Ronald N. Zuckermann¹ & Stephen Whitelam¹

A promising route to the synthesis of protein-mimetic materials that are capable of complex functions, such as molecular recognition and catalysis, is provided by sequence-defined peptoid polymers^{1,2}—structural relatives of biologically occurring polypeptides. Peptoids, which are relatively non-toxic and resistant to degradation³, can fold into defined structures through a combination of sequence-dependent interactions^{3–8}. However, the range of possible structures that are accessible to peptoids and other biological mimetics is unknown, and our ability to design protein-like architectures from these polymer classes is limited⁹. Here we use molecular-dynamics simulations, together with scattering and microscopy data, to determine the atomic-resolution structure of the recently discovered peptoid nanosheet, an ordered supramolecular assembly that extends macroscopically in only two dimensions. Our simulations show that nanosheets are structurally and dynamically heterogeneous, can be formed only from peptoids of certain lengths, and are potentially porous to water and ions. Moreover, their formation is enabled by the peptoids' adoption of a secondary structure that is not seen in the natural world. This structure, a zigzag pattern that we call a

Σ ('sigma')-strand, results from the ability of adjacent backbone monomers to adopt opposed rotational states, thereby allowing the backbone to remain linear and untwisted. Linear backbones tiled in a brick-like way form an extended two-dimensional nanostructure, the Σ -sheet. The binary rotational-state motif of the Σ -strand is not seen in regular protein structures, which are usually built from one type of rotational state. We also show that the concept of building regular structures from multiple rotational states can be generalized beyond the peptoid nanosheet system.

The peptoid nanosheet is a recently discovered, free-floating planar assembly that is only two molecules thick but that extends laterally for micrometres (Supplementary Figs 1, 2)^{8,10,11}. Nanosheets assemble from peptoids bearing alternating aromatic and charged sidechains (Fig. 1a) via compression at an air–water interface^{8,10–12}. We find, through a combination of atomic-force microscopy (AFM; Supplementary Fig. 3)^{8,10,13} and powder X-ray diffraction (XRD)^{8,10,11,13}, that nanosheets are bilayers, 3.0 ± 0.3 nm (\pm s.d.) thick, in whose interior the aromatic sidechains are sequestered, and on whose surfaces the charged sidechains are presented. Optical microscopy shows nanosheets to be

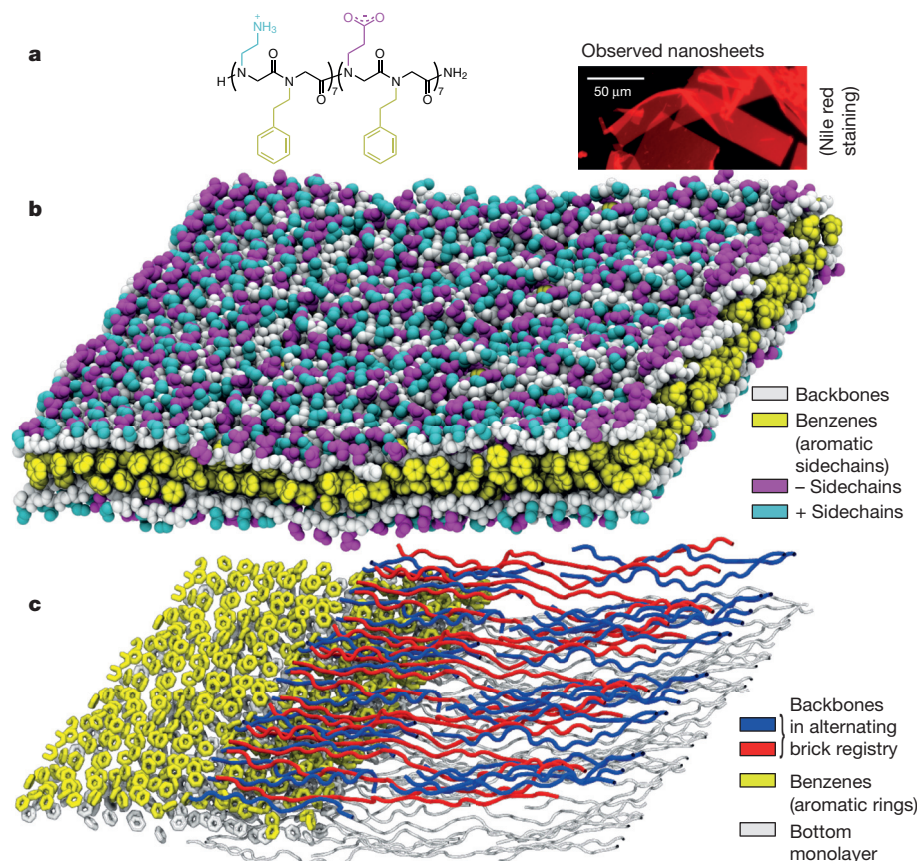


Figure 1 | Snapshot of a peptoid nanosheet obtained from molecular-dynamics simulations. **a**, Left, an amphiphilic 28-residue peptoid, which assembles into extended nanosheets only two molecules thick¹⁰, as shown in the fluorescent-microscopy image to the right. **b**, Snapshot of a bilayer, obtained from molecular-dynamics simulations. **c**, Backbones (blue and red) are generally collinear; aromatic rings (yellow) in the interior of the nanosheet show little order. (The lower leaf of the bilayer is coloured white.)

¹Molecular Foundry, Lawrence Berkeley National Laboratory, Berkeley, California 94709, USA. ²Center for Genomics and Systems Biology, New York University Abu Dhabi, Abu Dhabi, United Arab Emirates.

roughly rectangular, many square micrometres in horizontal extent, and to have relatively straight edges (Fig. 1a and Supplementary Fig. 2b). XRD measurements in the plane of the bilayer show peaks at characteristic distances of 4.5 Å and 3.6 Å (Supplementary Fig. 9), suggesting a high degree of order on the molecular scale. Direct observation of polymer chains in the nanosheet by transmission electron microscopy indicates a polymer–polymer parallel spacing of about 4.5 Å (ref. 8), in accordance with the 4.5-Å XRD peak.

To develop an atomistic model of peptoid nanosheets consistent with these observations, we used atomistic molecular-dynamics

simulations in conjunction with our recently developed CHARMM-based¹⁴ force field for peptoid backbones, MFTOID¹⁵ (CHARMM, Chemistry at HARvard Molecular Mechanics; MFTOID, Molecular Foundry (MF) and Peptoid (TOID)). We surveyed a range of low-energy nanosheet configurations (Supplementary Fig. 4a) as starting points for molecular-dynamics simulations; here, we report the results of one such set of simulations (see Methods).

Figure 1b and c show snapshots of a nanosheet patch, with periodic images displayed (for an image of the simulation box only, see Supplementary Fig. 4c). The nanosheet possesses linear order in its

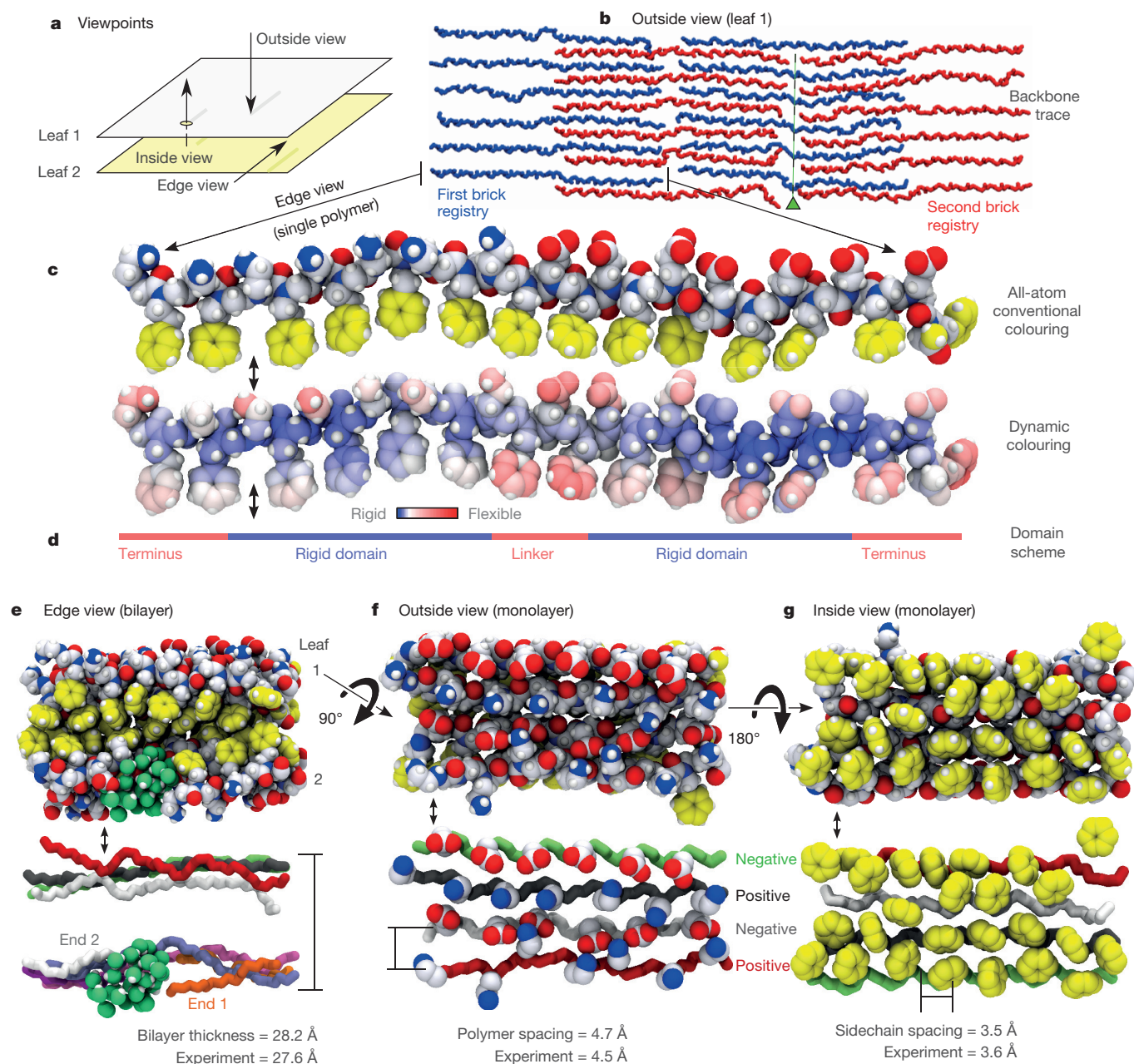


Figure 2 | Selected nanosheet features and dynamics. **a**, Three ways of viewing the bilayer: outside, inside and edge views. **b**, The general layout of polymers within each leaf is brick-like, with pockets (green triangle) formed at peptoid termini. **c**, **d**, Each peptoid, such as that shown here, can be coloured with conventional atomic colours (**c**, top), or with respect to dynamism (**c**, bottom; expanded in Supplementary Fig. 12), as measured by root-mean-squared fluctuations of atomic positions, showing that backbones possess three dynamically distinct regions (**d**): a rigid domain, a dynamic central linker and dynamic termini. As in globular proteins, the backbone is most dynamic in

termini and linker regions. **e–g**, Sections of the nanosheet from various angles and in various representations (enlarged in Supplementary Fig. 13). Each top and bottom panel (related by a double-headed arrow) show the same section and viewpoint, rendered in different ways. **e**, Water (green) encroaches at a pocket, suggesting that water channels might be engineered into nanosheets (Supplementary Fig. 20). **f**, Backbone order. **g**, Aromatic disorder. Real-space measurements taken from simulation averages and inferred from scattering experiments (shown at the bottom) agree to within 1 Å.

polymer backbones, and relative disorder in the aromatic and charged sidechains that are, respectively, internal and external to the bilayer (Supplementary Figs 11 and 12). Individual peptoid chains maintain the brick-like arrangement that is favoured by the segregation of charged sidechains into positively and negatively charged blocks (Fig. 1a), in which gaps or ‘pockets’ between polymer termini alternate with the backbones’ central portions, as one looks along the direction perpendicular to that of the polymer chains (Fig. 2b). As shown in Fig. 2c and Supplementary Fig. 12, backbone regions away from the pockets are the least dynamic portions of each peptoid, displaying a root-mean-squared fluctuation (RMSF) of ≤ 1 Å. The pocket-forming termini of each peptoid are more flexible than these interior regions, and are visibly less ordered (Fig. 2c, d and Supplementary Fig. 12), while the polymer sidechains are more flexible still. This hierarchy of flexibility is similar to that seen in proteins^{16,17}.

Figure 2e–g and Supplementary Fig. 8 show that the nanosheet thickness, interpolymer spacings and intersidechain spacings seen in our simulations lie within about an ångström of the mean characteristic distances seen in AFM and X-ray-scattering experiments^{8,10,11}.

The range of nanosheet thicknesses measured by AFM, 2.7 nm to 3.3 nm (Supplementary Fig. 3), is also reproduced in our simulations (Supplementary Fig. 8). These comparisons allow us to assign physical features to experimental measurements, and to verify the accuracy of our simulations (see Supplementary Table 1).

This brick-like arrangement of backbones that allows polymers to alternate with pockets has important consequences for the stability of nanosheets as a function of peptoid length. As polymer length decreases, it should become less energetically favourable for each residue to form a nanosheet, because pockets—near which peptoids possess fewer favourable electrostatic and aromatic interactions—increase in number per unit area as polymer length decreases. Simulations indeed show that short polymers are not stable in nanosheet form. In Fig. 3, we report the results of simulations done using nanosheets built from polymers that are 4, 8, 12, 16, 20, 24 or 28 residues long. Polymers longer than 12 residues form stable nanosheets. Nanosheets built from 12-residue polymers show signs of instability and a decrease in order (Fig. 3a–c) upon simulation, and nanosheets built from polymers shorter than 12 residues display an almost complete loss of structure within 40 ns of the start of the simu-

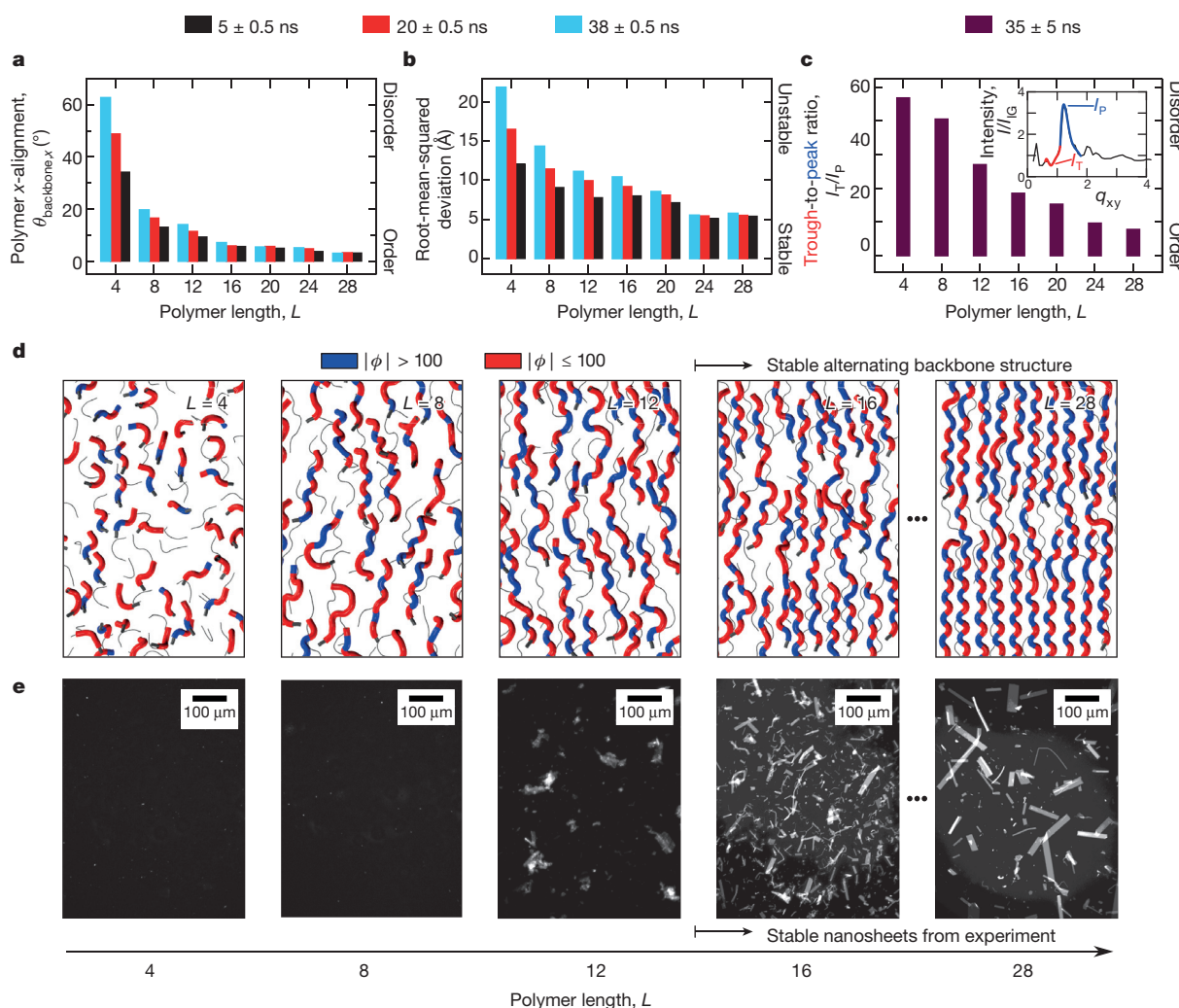


Figure 3 | Peptoids shorter than 12 residues do not form stable nanosheets. **a–c**, Metrics of stability and order for nanosheets built from peptoids of various lengths, L . These metrics are: **a**, mean polymer alignment with the x axis (equation (2) in Methods; see Supplementary Fig. 16a); **b**, root-mean-squared deviation from the initial configuration (heavy atoms only; see Supplementary Fig. 16b); and **c**, trough-to-peak ratio associated with backbone–backbone order obtained from simulated X-ray scattering (a spectrum for $L = 28$ is shown in the inset; see Supplementary Fig. 17). Larger values of all metrics indicate a lesser degree of order. Short-polymer nanosheets continue to evolve throughout the times simulated (**a**, **b**). I , intensity; I_p , peak intensity; I_t , trough

intensity; I_{IG} , intensity expected of an ideal gas (for the purposes of normalization); q_{xy} , xy component of the wavevector. **d**, **e**, Simulations (**d**) and experiments (fluorescence optical microscopy; **e**) done with peptoids of lengths $L = 4, 8, 12, 16$ or 28 show a similar trend: peptoids shorter than 12 residues do not form stable nanosheets. The case $L = 12$ is marginal: in experiments, high concentrations of peptoids result in the formation of small nanosheets, which become disordered on a timescale of about one week (see Supplementary Fig. 18). Colouring backbones with respect to the backbone dihedral angle, ϕ , makes clear that adjacent backbone monomers adopt opposed rotational states. This novel secondary structure, the Σ -strand, is discussed in Fig. 4.

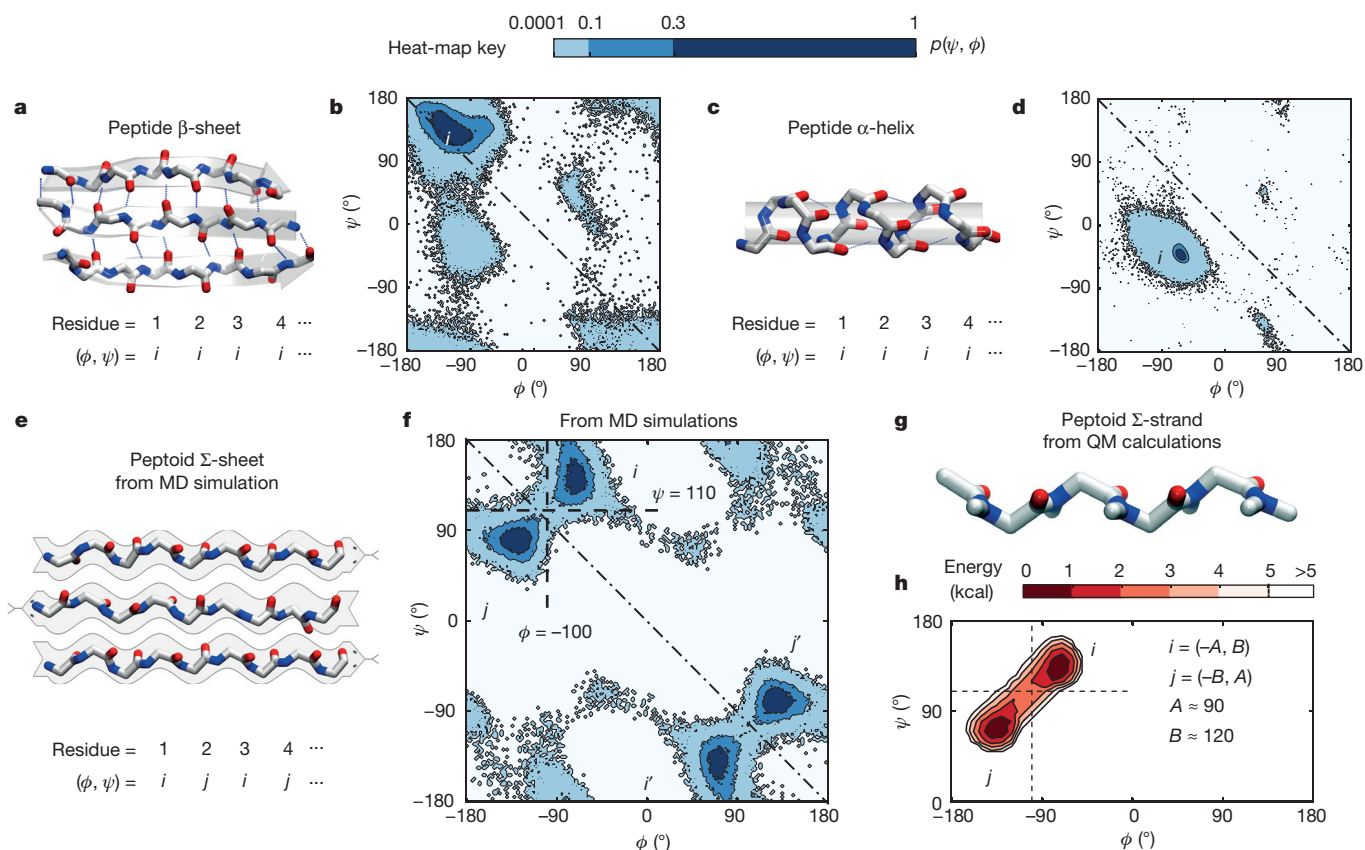


Figure 4 | Peptoid nanosheets are made possible by a novel secondary structure. **a–d**, Typical secondary structures found in proteins, such as the β -sheet (**a**) and α -helix (**c**), are described by specific pairs of backbone dihedral angles, ϕ and ψ , whose position is labelled by i in their respective Ramachandran plots (**b**, **d**; see Methods). The heat map is a two-dimensional histogram showing the normalized probability $p(\psi, \phi)$ that a residue in one of the indicated structures will adopt particular combinations of the dihedral angles ϕ and ψ . **e**, **f**, In contrast, the secondary structure observed in the peptoid nanosheet (**e**) consists of two characteristic (ϕ, ψ) positions, i and j (**f**). (The other occupied regions in **f**, that is, the pale blue locations, indicate polymer ends and linker regions; see Supplementary Figs 21 and 22. Owing to the

lateral (Fig. 3d and Supplementary Fig. 15). In accordance with our simulations, our experiments show that peptoids that are 12 or fewer residues in length do not form stable nanosheets (Fig. 3e and Supplementary Figs 18 and 19). Our simulations suggest that this failure occurs because short-polymer nanosheets are unstable relative to disordered aggregates. Our simulations also suggest that if polymer connections are reinforced—for example, by head-to-tail cross-linking—this limitation could be overcome.

When considering the design of nanosheets, the observation that interchain pockets are tolerated by polymers longer than a certain limit shows that pockets should be considered to be an integral part of the nanosheet structure. Simulations show that such pockets allow water to encroach on the aromatic centre of the bilayer (Fig. 2e and Supplementary Fig. 20b). This observation suggests that peptoid termini could be modified to create pockets able to bind specific small molecules, potentially permitting catalytic function. Furthermore, our simulations show that, if pockets on opposing leaves coincide, then nanosheets possess channels through which water can pass (Supplementary Fig. 20c), indicating the potential of nanosheets as selective membranes. Experimental work to test these predictions is under way.

The stability of extended, planar peptoid nanosheets is enabled by the polymers' linear, untwisted configurations. Our simulations reveal that this linearity results from the ability of sequential backbone residues to adopt one of two states, whose rotation about the backbone axis opposes and cancels each other. In Fig. 3d we have coloured these two

achirality of the peptoid backbone's α -carbon, i' and j' are equivalent to i and j , respectively.) Positions i and j are visibly equidistant from the achiral diagonal of the Ramachandran plot, and so adopting the two states in an alternating fashion allows the peptoid backbone to remain linear and untwisted. This motif, the Σ -strand, allows the formation of extended planar structures that have not been made previously using protein-like building motifs. **g**, Quantum-mechanical (QM) calculations of a range of possible Σ -type strands show energy minima (**h**) that match the high-occupancy regions of **f**, indicating that the rotational tendencies of isolated polymers are preserved by the sidechain interactions established within the nanosheet. (MD, molecular dynamics. The alternating pattern is defined by the numbers A and B .)

rotational states red and blue. This building principle is distinct from that used by proteins: the backbones of protein secondary structures such as α -helices and β -sheets are defined primarily by a single rotational state.

Protein rotational states are quantified by their backbone dihedral angles, traditionally denoted ϕ and ψ , and conventionally described by a Ramachandran plot¹⁸. As shown in Fig. 4a–d, regular chiral protein structures such as the α -helix and β -sheet correspond roughly to a single location, i , on the Ramachandran plot¹⁸. In contrast, stable nanosheets are composed of peptoids whose backbone states occupy two specific regions of the Ramachandran plot (Fig. 4f), labelled i and j . Chains whose adjacent residues alternate between these two states remain linear. A snapshot of three backbone segments (Fig. 4e) emphasizes this alternating motif. We call this motif the Σ -strand, because its linear, twist-free nature derives from the combination, or sum (Σ), of its two rotational states (and because the resulting polymer 'snakes' back and forth). In principle one could have a Σ -strand built from any two opposed rotational states. However, density functional theory calculations show that the particular rotational states observed in our atomistic simulations are the lowest-energy Σ -type arrangement for isolated polymers (Fig. 4g and Supplementary Figs 23–26). This comparison provides additional confidence in the accuracy of the MFTOID force-field simulations, and confirms that the basic rotational tendency of an isolated peptoid backbone¹⁵ (Fig. 4h) is preserved by the side-chain–sidechain interactions established within the nanosheet (Fig. 4f).

Single-point secondary structures

Higher-order (two-point) secondary structures

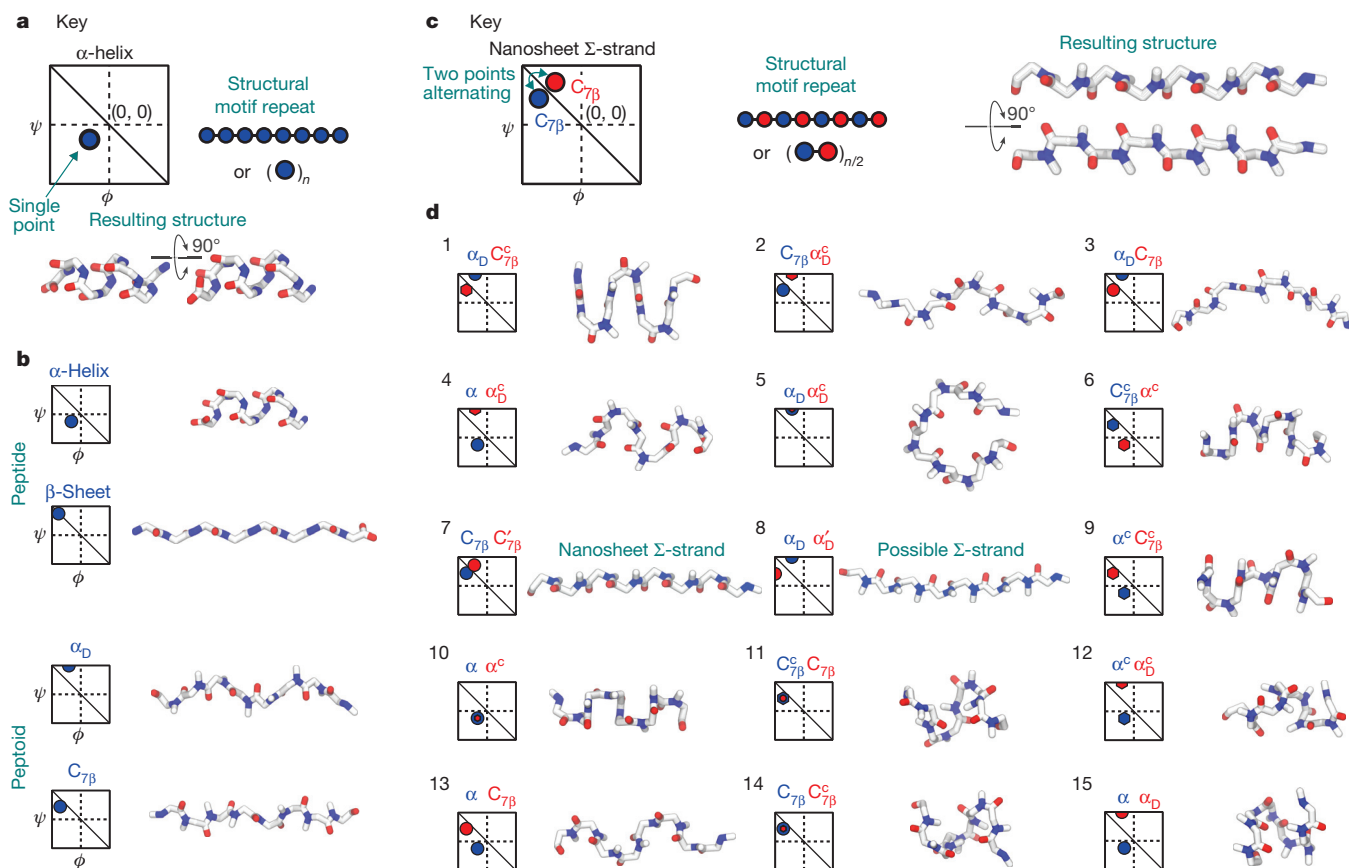


Figure 5 | The design principle underlying the peptoid nanosheet has general application. **a, b,** Canonical secondary structures are ‘single-point’ ones, meaning that they are described by single regions on a Ramachandran plot. ‘D’ and ‘7 β ’ are labels given to certain common peptoid structures. **c,** Here we propose the concept of building higher-order or multi-point secondary structures using combinations of points (alternating blue and red dots) on the Ramachandran plot. **d,** This principle can be used to create a large set of novel structures, of which the Σ -strand is but one example, potentially expanding the range of functional polymer structures. For instance, model 8 shows that a linear chain can display alternating orthogonal sidechains, which may allow the

formation of a planar assembly distinct from the Σ -sheet (‘possible Σ -strand’). Model 5, if cyclized, may allow for the assembly of stacked disks (alternating rotational properties have been observed in cyclic peptoids²², while alternating L,D cyclic peptides have been found to stack²³). As in the case of the Σ -strand (model 7), such structures can in principle be stabilized by physically reasonable sidechain interactions possessing appropriate geometries. Supplementary Table 2 discusses the method by which the structures in **d** were built, relaxed and ranked. For each key in **d**, the superscript ‘c’ indicates a *cis*-amide-bond dihedral angle.

The Σ -strand motif is protein-like in its regularity, but is different from protein secondary structures in important ways: it is a molecular motif that permits planar assemblies of macroscopic extent; it is stabilized for reasons other than hydrogen bonding; and it is built from two (not one) rotational states. The closest comparable secondary structures in proteins, β -sheets, are unable to maintain macroscopic flatness owing to a propeller-like twist in the shape of each strand (although fibril-like protein assemblies can be of considerable size¹⁹ and can be combined to form structures of macroscopic extent²⁰). The Σ -sheet therefore represents a complement to the basic building blocks of the protein world. To replicate the Σ -sheet assembly using other polymers, one therefore needs two things: molecular rotational properties that give rise to symmetry about the diagonal of the Ramachandran plot (Fig. 4f), and sidechain sequence patterning that promotes these conformations within a desired assembly.

The principle underlying the stability of the Σ -strand and the peptoid nanosheet is that regular secondary structures can be built from more than one rotational state—that is, from more than one point on the Ramachandran plot. This principle immediately suggests the possibility of a large set of novel secondary structures, some of which are shown in Fig. 5. The potential applications of this set of structures go beyond those of the nanosheet system—the Σ -strand is but one member of this set—and may include the formation of novel folded structures and

assemblies. Since the 1950s, the literature on protein secondary structure has focused on the idea of building with one type of rotational state¹⁸, to make ‘single-point’ Ramachandran structures. We propose that building higher-order, multi-point secondary structures might greatly expand the repertoire of folded polymeric building blocks. Figure 5d provides guidelines for backbone design. It may also help to characterize structurally the rapidly growing family of solid-state peptoid polymer crystals²¹. Such crystals consist largely of extended polymer conformations that so far have eluded characterization at atomic resolution; the extended backbones seen in Fig. 5d are candidate conformations. More generally, the new building principle that we have identified, combined with the ability to encode into peptoids a defined sequence of chemically diverse monomers, offers a way to create new structured polymers through combinatorial design.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 20 April 2015; accepted 27 July 2015.

Published online 7 October 2015.

1. Knight, A. S., Zhou, E. Y., Francis, M. B. & Zuckermann, R. N. Sequence programmable peptoid polymers for diverse materials applications. *Adv. Mater.* <http://dx.doi.org/10.1002/adma.201500275> (2015).

2. Zuckermann, R. N. & Kodadek, T. Peptoids as potential therapeutics. *Curr. Opin. Mol. Ther.* **11**, 299–307 (2009).
3. Sun, J. & Zuckermann, R. N. Peptoid polymers: a highly designable bioinspired material. *ACS Nano* **7**, 4715–4732 (2013).
4. Butterfoss, G. L. *et al.* De novo structure prediction and experimental characterization of folded peptoid oligomers. *Proc. Natl Acad. Sci. USA* **109**, 14320–14325 (2012).
5. Yoo, B. & Kirshenbaum, K. Peptoid architectures: elaboration, actuation, and application. *Curr. Opin. Chem. Biol.* **12**, 714–721 (2008).
6. Laursen, J. S., Engel-Andreasen, J., Fristrup, P., Harris, P. & Olsen, C. A. *Cis-trans* amide bond rotamers in β -peptoids and peptoids: evaluation of stereoelectronic effects in backbone and side chains. *J. Am. Chem. Soc.* **135**, 2835–2844 (2013).
7. Butterfoss, G. L., Renfrew, P. D., Kuhlman, B., Kirshenbaum, K. & Bonneau, R. A preliminary survey of the peptoid folding landscape. *J. Am. Chem. Soc.* **131**, 16798–16807 (2009).
8. Nam, K. T. *et al.* Free-floating ultrathin two-dimensional crystals from sequence-specific peptoid polymers. *Nature Mater.* **9**, 454–460 (2010).
9. Drexler, K. E. Peptoids at the 7th Summit: toward macromolecular systems engineering. *Biopolymers* **96**, 537–544 (2011).
10. Sanii, B. *et al.* Shaken, not stirred: collapsing a peptoid monolayer to produce free-floating, stable nanosheets. *J. Am. Chem. Soc.* **133**, 20808–20815 (2011).
11. Kudirka, R. *et al.* Folding of a single-chain, information-rich polypeptoid sequence into a highly ordered nanosheet. *Biopolymers* **96**, 586–595 (2011).
12. Sanii, B. *et al.* Structure-determining step in the hierarchical assembly of peptoid nanosheets. *ACS Nano* **8**, 11674–11684 (2014).
13. Robertson, E. J. *et al.* Assembly and molecular order of two-dimensional peptoid nanosheets through the oil-water interface. *Proc. Natl Acad. Sci. USA* **111**, 13284–13289 (2014).
14. Brooks, B. R. *et al.* CHARMM: the biomolecular simulation program. *J. Comput. Chem.* **30**, 1545–1614 (2009).
15. Mirijanian, D. T., Mannige, R. V., Zuckermann, R. N. & Whitelam, S. Development and use of an atomistic CHARMM-based forcefield for peptoid simulation. *J. Comput. Chem.* **35**, 360–370 (2014).
16. Karplus, M. & McCammon, J. A. The internal dynamics of globular proteins. *CRC Crit. Rev. Biochem.* **9**, 293–349 (1981).
17. Halle, B. Flexibility and packing in proteins. *Proc. Natl Acad. Sci. USA* **99**, 1274–1279 (2002).
18. Berg, J. M., Tymoczko, J. L. & Stryer, L. *Biochemistry, International Edition* 7th edn (WH Freeman & Co., 2010).
19. Hammer, N. D., Wang, X., McGuffie, B. A. & Chapman, M. R. Amyloids: friend or foe? *J. Alzheimers Dis.* **13**, 407–419 (2008).
20. Knowles, T. P., Oppenheim, T. W., Buell, A. K., Chirgadze, D. Y. & Welland, M. E. Nanostructured films from hierarchical self-assembly of amyloidogenic proteins. *Nature Nanotechnol.* **5**, 204–207 (2010).
21. Sun, J., Teran, A. A., Liao, X., Balsara, N. P. & Zuckermann, R. N. Crystallization in sequence-defined peptoid diblock copolymers induced by microphase separation. *J. Am. Chem. Soc.* **136**, 2070–2077 (2014).
22. Maulucci, N. *et al.* Synthesis, structures, and properties of nine-, twelve-, and eighteen-membered N-benzyloxyethyl cyclic α -peptoids. *Chem. Commun.* **33**, 3927–3929 (2008).
23. Zhang, W. *et al.* PEG-stabilized bilayer nanodisks as carriers for doxorubicin delivery. *Mol. Pharm.* **11**, 3279–3290 (2014).

Supplementary Information is available in the online version of the paper.

Acknowledgements Portions of this work were done as a User project at the Molecular Foundry at Lawrence Berkeley National Laboratory, supported by the Office of Science, Office of Basic Energy Sciences, of the US Department of Energy under contract no. DE-AC02-05CH11231. R.V.M., T.K.H., C.P., E.J.R., A.B., R.N.Z. and S.W. were supported by the Defense Threat Reduction Agency under contract no. IACRO-B0845281. C.P. was also supported by the Natural Sciences and Engineering Research Council of Canada (NSERC PDF). R.N.Z. and S.W. were also supported by the Office of Science, Office of Basic Energy Sciences, of the US Department of Energy under contract no. DE-AC02-05CH11231. We thank G. K. Olivier for providing the AFM data. This work used resources of the National Energy Research Scientific Computing Center, which is supported by the Office of Science of the US Department of Energy under contract no. DE-AC02-05CH11231. Quantum-mechanical calculations were carried out on the High Performance Computing resources at New York University Abu Dhabi.

Author Contributions R.V.M., R.N.Z. and S.W. initiated the research. R.V.M. and S.W. designed the molecular-dynamics simulations; R.V.M. performed the simulations. T.K.H. performed simulated X-ray-scattering calculations; C.P., E.J.R. and A.B. performed the experiments; G.L.B. designed and performed the quantum-mechanical calculations. All authors contributed to analysing the results and writing the paper.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to R.V.M. (rvmannige@lbl.gov) or S.W. (swwhitelam@lbl.gov).

METHODS

Experimental nanosheet synthesis. Block-charge peptoids of lengths $L = 4, 8, 12, 16$ or 28 residues were synthesized and purified using automated solid-phase synthesis³. For $L \geq 12$ residues (Supplementary Fig. 18d–f) and $L > 12$ residues (Fig. 3e), peptoids were present at a concentration of 20 μM in 10 mM Tris buffer, pH 8.0. Peptoid nanosheets of length $L \leq 12$ were absent at this concentration (Supplementary Fig. 18d), perhaps because not enough polymer adsorbed to the air–water interface (see discussions following Supplementary Fig. 18d); therefore we attempted to produce $L \leq 12$ nanosheets at higher concentrations (Fig. 3e and Supplementary Fig. 18a–c). Peptoids of length $L = 12$ were present at a concentration of 1 mM (Supplementary Fig. 18c), of length $L = 8$ at 22 mM (Supplementary Fig. 18b), and of length $L = 4$ at 41 mM (Supplementary Fig. 18a), all in 10 mM Tris buffer, pH 8.0. Sheets were prepared by agitation via the vial-rocking method¹⁰. To test whether the peptoids were adsorbing to the air–water interface—a prerequisite to assembly—we obtained surface tension data for the $L = 4, 8$ and 12 peptoids via the pendant drop method (see discussions following Supplementary Fig. 18)²⁴.

Nanosheet yield analysis by optical microscopy. Identical volumes were removed from each sample and applied to a thin agarose-gel slice to allow imaging for Fig. 3e and Supplementary Fig. 18 (ref. 25).

Measuring distances by single-angle X-ray scattering and AFM. Nanosheet thicknesses were calculated via AFM^{8,10,13}. We obtained AFM images of dry nanosheets deposited on a mica substrate in ambient air. Other distances, such as intrapolymer spacings (Fig. 2 and Supplementary Fig. 8), were obtained via XRD^{8,10,11,13}. XRD data also show a peak corresponding to the thickness of the nanosheet.

Molecular-dynamics protocol. We surveyed a range of low-energy nanosheet configurations (Supplementary Fig. 4a) as starting points for molecular-dynamics simulations. To conduct this survey we built 28-residue polymers into monolayers in a brick-like arrangement, and joined two identical opposing monolayers to form a bilayer (Supplementary Fig. 4a). The brick arrangement is suggested by simple electrostatic considerations, because peptoids' charged sidechains are segregated into positive and negative blocks along each chain¹¹ (Fig. 1a). The computational search space so defined possesses six degrees of freedom, described by distances between polymers within the same layer and between monolayers (Supplementary Fig. 4a). We calculated energies (in implicit water) for a large number of these parameter sets, and chose six low-energy versions on which to focus (Supplementary Fig. 4b). Each nanosheet version was then explicitly solvated and relaxed in a sequence of protocols (see Supplementary Fig. 4d), before undergoing constant-pressure molecular-dynamics simulation using a leap-frog integrator (a commonly used second-order numerical method for integrating equations of motion) with a 1-fs timestep.

Each of the six low-energy versions of the nanosheet was seen to be stable for over 50 ns of molecular-dynamics simulation at standard temperature and pressure; that is, each had ceased to evolve (structural snapshots for each simulation are available in Supplementary Information), and did not dissolve or convert to a different structure (Supplementary Fig. 5). (In contrast, high-energy initial configurations did not yield stable structures.) These simulations reveal that, in the region of the energy minimum dictated by charged sidechains¹¹, there exists a rugged free-energy landscape with a range of near-degenerate nanosheet structures. Biomolecules, by contrast, often possess clear free-energy minima^{26–30}.

Our simulations represent periodically replicated patches of nanosheets of approximate dimensions 60 nm \times 18 nm, and so do not address the nature of nanosheet order on the micrometre scale. Given that nanosheets are produced by a far-from-equilibrium mechanical protocol¹⁰, and given the stability of all six nanosheet versions in our simulations, it is possible that extended nanosheets consist of a patchwork of different types of stable local order. Nonetheless, each of the six low-energy nanosheet versions displays molecular features consistent with experiments (Supplementary Fig. 8), from which we infer that substantial portions of nanosheets display locally the atomic-scale features seen in our simulations. In the main text, for brevity, we present results from one particular nanosheet version, number 5.

Supplementary Figure 4 describes our generation of the six molecular-dynamics starting configurations. Molecular-dynamics simulations were done using the CHARMM software package¹⁴. Aside from the initial set-up, all nanosheet series were simulated using a leap-frog integration algorithm with a 1-fs timestep. Simulations were performed in the isothermal-isobaric (NPT, for constant particle number, pressure and temperature) ensemble at 300 K and 1 atmospheric pressure, in an orthorhombic periodic box in which the three orthogonal box dimensions were allowed to vary independently. The Hoover algorithm was used to maintain constant pressure³¹. Hydrogen bonds were constrained using the SHAKE protocol³². Particle-mesh Ewald summation was used to evaluate long-range electrostatic interactions, with a real-space cut-off of 12 Å, a sixth-order

cubic spline, and a κ value (width of the Gaussian distribution) of 0.34. Van der Waals interactions were calculated up to a distance of 12 Å, with a smoothing function applied from 10 Å to 12 Å.

Simulations reported in the main text were performed in a solution consisting of water, the nanosheet's counterions, and potassium chloride at a concentration of 10 mM. We also carried out a series of independent simulations of the version 5 nanosheet, in which potassium chloride concentrations were set to 0.1 mM, 1 mM, 50 mM and 100 mM (Supplementary Fig. 14). These simulations verified that simulated nanosheet structures do not change markedly over the range of salt concentrations used in our experiments^{8,10,11}.

Computational X-ray scattering. We calculated scattering spectra at wavevector q using the expression

$$I(q) = \left| \sum_j f_j \exp(iq \cdot r_j) \right|^2 \quad (1)$$

where the sum runs over all atoms, f_j is the atomic scattering factor for atom j and r_j is the position of atom j . Equation (1) assumes that electrons are localized at atomic sites. We let f_j equal the atomic number, an approximation that is nearly exact at the experimental X-ray energies of 11 keV (ref. 33). The values of q were discretized so as to be commensurate with the periodic box.

Experimental in-plane X-ray-scattering spectra were taken by allowing nanosheets to dry on a Kapton grid, stacking the sheets on top of each other to produce a sample with a uniform orientation¹². The absence of lamellar peaks when the X-ray beam was fired 'face-on' into the stack confirmed that the nanosheet normals were uniformly oriented parallel to the grid normal. However, the radially averaged signal indicated that various sheets and/or domains existed within the beamline, as expected given that the beam cross-section (120 $\mu\text{m} \times$ 800 μm) is greater than the typical nanosheet size (20 $\mu\text{m} \times$ 20 μm). To compare with these in-plane, radially averaged spectra, we radially averaged the simulated scattering spectra (equation (1)) in the xy plane of the nanosheet, $I(q_{xy}) = \langle I(q) \rangle_{q_z=0, \sqrt{q_x^2 + q_y^2} = q_{xy}}$. Finally, these intensities ($I(q_{xy})$) were normalized by the expected intensity of an equivalent ideal gas ($I_{\text{IG}}(q_{xy})$), which is discussed in Supplementary Fig. 17.

Direct measurements. While X-ray scattering is useful for comparing simulation to experiment, direct measurements of some features (for example, polymer y -padding, nearest neighbour N–N spacing, bilayer thickness) provide additional information with regards to distances and distributions. Supplementary Fig. 7 and its associated text discusses those measurements and how they are calculated, while Supplementary Fig. 8 compares experimental observations to those metrics.

Measurement of order in the bilayer. For each peptoid (Fig. 3a), we designate the position of the backbone nitrogen of the i th residue to be r_i . A simple measure of order is the angle $\theta_{\text{backbone},x}$ between each peptoid backbone and the x axis

$$\theta_{\text{backbone},x} = \cos^{-1} \left[\frac{(r_{\text{last}} - r_{\text{second}}) \cdot \hat{x}}{|r_{\text{last}} - r_{\text{second}}|} \right] \quad (2)$$

α -Helical and β -sheet segments from proteins. A structure database was obtained from the Structural Classification of Proteins (SCOPe; release 2.03) that contains proteins with no more than 40% sequence identity with each other (<http://scop.berkeley.edu/downloads/pdbstyle/pdbstyle-sel-gs-bib-40-2.03.tgz>). α -Helical and β -sheet segments were identified using the DSSP (Define Secondary Structure of Proteins) algorithm³⁴. (ϕ, ψ) pairs were obtained for all residues within α -helical and β -sheet segments, which contributed towards the histograms in Fig. 4b, d, respectively.

Quantum mechanics calculations. Density functional theory was used to assess the low-energy structures available to a peptoid backbone consisting of two alternating sets of dihedral angles (ϕ, ψ) . B3LYP and M05-2X were the two functionals used to calculate the energies of a configuration. B3LYP is a commonly used and widely applicable functional³⁵, and M05-2X is a newer functional that accounts for dispersion forces²². Gaussian 09 (C.01) was used for density functional theory calculations²³. For each peptoid model, the alternating pattern is defined by the numbers A and B (Fig. 4h), where adjacent (ϕ, ψ) pairs take the values $(-A, B)$ and $(-B, A)$. A and B were constrained to range from 50° to 180° in steps of 10°. Molecules were optimized at the HF/6-31G* level of theory³⁶, with A and B constrained and all other degrees of freedom unconstrained. The lowest energy associated with each (A, B) pair is reported in the Ramachandran plots of Fig. 4h and Supplementary Figs 24–26.

24. Rosen, M. J. & Kunjappu, J. T. *Surfactants and Interfacial Phenomena* (John Wiley & Sons, 2012).

25. Leopold, P. E., Montal, M. & Onuchic, J. N. Protein folding funnels: a kinetic approach to the sequence-structure relationship. *Proc. Natl Acad. Sci. USA* **89**, 8721–8725 (1992).

26. Bryngelson, J. D., Onuchic, J. N., Socci, N. D. & Wolynes, P. G. Funnels, pathways, and the energy landscape of protein folding: a synthesis. *Proteins* **21**, 167–195 (1995).
27. Schug, A. & Onuchic, J. N. From protein folding to protein function and biomolecular binding by energy landscape theory. *Curr. Opin. Pharmacol.* **10**, 709–714 (2010).
28. Shan, Y. *et al.* How does a drug molecule find its target binding site? *J. Am. Chem. Soc.* **133**, 9181–9183 (2011).
29. Hoover, W. G. Canonical dynamics: equilibrium phase-space distributions. *Phys. Rev. A* **31**, 1695–1697 (1985).
30. Allen, M. P. & Tildesley, D. J. *Computer Simulation of Liquids* (Oxford Univ. Press, 1989).
31. Chantler, C. Theoretical form factor, attenuation, and scattering tabulation for $Z=1-92$ from $E=1-10$ eV to $E=0.4-1.0$ MeV. *J. Phys. Chem. Ref. Data* **24**, 71 (1995).
32. Kabsch, W. & Sander, C. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* **22**, 2577–2637 (1983).
33. Lee, C., Yang, W. & Parr, R. G. Development of the Colle-Salvetti correlation-energy formula into a functional of the electron density. *Phys. Rev. B* **37**, 785–789 (1988).
34. Zhao, Y., Schultz, N. E. & Truhlar, D. G. Exchange-correlation functional with broad accuracy for metallic and nonmetallic compounds, kinetics, and noncovalent interactions. *J. Chem. Phys.* **123**, 161103 (2005).
35. Frisch, M. J. *et al.* *Gaussian09* Revision D.01 (Gaussian Inc., Wallingford, 2009).
36. Roothaan, C. C. J. New developments in molecular orbital theory. *Rev. Mod. Phys.* **23**, 69–99 (1951).

The multi-millennial Antarctic commitment to future sea-level rise

N. R. Golledge^{1,2}, D. E. Kowalewski³, T. R. Naish^{1,2}, R. H. Levy², C. J. Fogwill⁴ & E. G. W. Gasson⁵

Atmospheric warming is projected to increase global mean surface temperatures by 0.3 to 4.8 degrees Celsius above pre-industrial values by the end of this century¹. If anthropogenic emissions continue unchecked, the warming increase may reach 8–10 degrees Celsius by 2300 (ref. 2). The contribution that large ice sheets will make to sea-level rise under such warming scenarios is difficult to quantify because the equilibrium-response timescale of ice sheets is longer than those of the atmosphere or ocean. Here we use a coupled ice-sheet/ice-shelf model to show that if atmospheric warming exceeds 1.5 to 2 degrees Celsius above present, collapse of the major Antarctic ice shelves triggers a centennial- to millennial-scale response of the Antarctic ice sheet in which enhanced viscous flow produces a long-term commitment (an unstoppable contribution) to sea-level rise. Our simulations represent the response of the present-day Antarctic ice-sheet system to the oceanic and climatic changes of four representative concentration pathways (RCPs) from the Fifth Assessment Report of the Intergovernmental Panel on Climate Change³. We find that substantial Antarctic ice loss can be prevented only by limiting greenhouse gas emissions to RCP 2.6 levels. Higher-emissions scenarios lead to ice loss from Antarctic that will raise sea level by 0.6–3 metres by the year 2300. Our results imply that greenhouse gas emissions in the next few decades will strongly influence the long-term contribution of the Antarctic ice sheet to global sea level.

Ice sheets lose mass and contribute to global mean sea level (GMSL) by surface and basal melting, and by dynamic thinning^{4,5}. Melting typically occurs at the ice surface if air temperatures are above zero, and at the ice base either where grounded ice is at the pressure melting point, or where ice is afloat in the ocean⁶. Owing to the exchange of heat between ice and adjacent water or air, changes in environmental temperature may cause changes in melt rate. Dynamic ice loss involves large-scale adjustment of part of an ice sheet to a change in the force balance that determines the ice flow speed (for example, loss of a buttressing ice shelf). There may be a delay in the dynamic ice-sheet response because although perturbations may be transmitted rapidly through ice shelves, the grounded ice sheet takes time to relax to a new steady state. Understanding which environmental changes lead to an immediate, and which to a lagged, ice-sheet response is important in predicting the timescale and transient evolution of the perturbed ice-sheet contribution to GMSL, as well as the total contribution that anthropogenic greenhouse-gas emissions will ultimately commit the planet to.

In Antarctica, peripheral ice shelves restrain the flow of grounded ice⁷ and are sensitive to warming air and ocean temperatures^{8–10}. Observations show that 67% to 98% of ocean warming since 2006 occurred in the Southern Ocean¹¹, with a similar multidecadal trend in the heat content of circum-Antarctic waters¹². Air temperatures across Antarctica, especially in West Antarctica, have risen on average 0.6 °C in the past 50 years¹³. Precipitation has also increased, but there remains a net acceleration of mass loss from Antarctica¹⁴ that continues to contribute to rising sea levels¹⁵. Recent studies indicate

that parts of West Antarctica may already be undergoing irreversible retreat^{16,17}, and several studies^{18,19} have focused on quantifying the resultant Antarctic contribution to GMSL by 2100. Predicting the multi-centennial-scale to millennial-scale commitments implied by future climate scenarios²⁰ and palaeoclimate equilibrium ice-sheet reconstructions for high-carbon-dioxide (CO₂) conditions^{21,22} has received less attention, however.

Some studies have investigated long-term commitments to sea-level rise under global warming scenarios using statistical relationships between past temperatures and global sea levels^{20,23}. Here we present numerical simulations of the Antarctic ice-sheet/ice-shelf system response to environmental changes predicted by four RCPs of the Fifth Assessment Report of the Intergovernmental Panel on Climate Change³ (Extended Data Fig. 1; Extended Data Table 1), which we extend to 5000 CE to capture the multi-millennial response (see 'Methods' and Extended Data Fig. 2). We use the Parallel Ice-Sheet Model, an open-source, three-dimensional, thermodynamic, coupled ice-sheet/ice-shelf model^{24,25}, and run simulations at spatial resolutions of 10 km and 20 km. Our method employs two different grounding-line parameterizations to quantify the likely range of ice-sheet responses. One implementation uses sub-grid interpolation of basal melting at grounding lines²⁶ whereas the other does not²⁷ (see Methods). Because the former tends to accelerate grounding-line retreat in coarse-resolution models such as ours we refer to sea-level contributions from these simulations as 'high' and those that do not include the sub-grid basal melt interpolation as 'low'. Other schemes may produce even higher or lower values, however.

The most immediate response to predicted climate changes in all but our RCP 2.6 experiment is a reduction in extent of the major ice shelves (Ross, Filchner-Ronne and Amery) within 100–300 years (Fig. 1). Reduced ice-shelf buttressing leads to increased discharge and flow acceleration that promotes grounding-line recession in areas where marine-based ice sheets occupy deep basins (such as West Antarctica; Fig. 1e, h, k). In our 'high' simulations, prolonged warming also leads to substantial loss of ice in East Antarctica, particularly in the Wilkes subglacial basin, at the margin of the Aurora subglacial basin, and in the southern and eastern Weddell Sea embayment (pale blue shading and blue lines in Fig. 1f, i and l). Under RCP 8.5 the range of sea-level contributions predicted by our 'low' and 'high' simulations respectively is 0.1–0.39 m by 2100, increasing to 1.6–2.96 m by 2300 and 5.2–9.31 m by 5000 CE (Figs 1j–l and 2a). Rates of sea-level rise are 5.5–15 mm per year by 2300 under RCP 8.5 conditions, but even under the lesser forcings of RCP 4.5 and RCP 6 reach 3–5 mm per year by 2300 (Figs 1e, h and 2a). These simulations show that systemic lags delay the fastest rates of ice-sheet loss by decades to centuries following the onset of initial forcing (Fig. 2a).

The initial ice-sheet/ice-shelf responses are important, but are relatively modest compared to the responses triggered now that will take place (committed changes) over subsequent millennia (Fig. 1c, f, i, l). In our experiments we see clear correspondence between the magnitude of

¹Antarctic Research Centre, Victoria University of Wellington, Wellington 6140, New Zealand. ²GNS Science, Avalon, Lower Hutt 5011, New Zealand. ³Department of Earth, Environment, and Physics, Worcester State University, Worcester, Massachusetts 01602, USA. ⁴Climate Change Research Centre, University of New South Wales, Sydney, New South Wales 2052, Australia. ⁵Climate System Research Center, University of Massachusetts Amherst, Amherst, Massachusetts 01003, USA.



Figure 1 | Modelled ice-sheet evolution under Antarctic-specific RCP-based warming scenarios. Emissions-forced climate warming to 2100 CE (a, d, g, j) and 2300 CE (b, e, h, k) results in initial sea-level contributions from Antarctica that are only a small proportion of the total sea-level commitment by 5000 CE (c, f, i, l). Magnitudes and rates of sea-level contributions are shown for each panel. Leading values and those in parentheses relate to 'low' and 'high' scenarios respectively. Ice extent for 'low' simulations is shown in

white; blue lines show grounding-line locations for 'high' simulations. Pale blue shading shows grounded ice lost in 'high' simulations but present in the 'low' scenario. Note the increasing divergence between 'high' and 'low' beyond 2300 CE. Grey texturing indicates areas of relatively faster-flowing ice. WAIS, West Antarctic Ice Sheet; EAIS, East Antarctic Ice Sheet; FRIS, Filchner–Ronne Ice Shelf.

a multi-millennial environmental perturbation, the steady-state area of fringing ice shelves, and the long-term ice-sheet contribution to GMSL (Fig. 2b). Specifically, we observe a sharp decline in near-equilibrium ice-shelf extent when atmospheric and oceanic temperatures are maintained 1.2 °C and 0.3 °C respectively above present. Almost all floating ice is lost if equilibrium air and ocean temperatures increase by more than 2 °C and 0.5 °C above present, respectively. These ice shelf 'thresholds' (defined here as an abrupt reduction in area to 50% of present) also occur during centennial-scale adjustment of the ice sheet to new equilibria, but happen at different temperatures (Fig. 2b). We infer therefore that short-lived (decadal-scale) environmental perturbations may mimic the longer-term responses if the former are of a large enough magnitude.

To isolate the causes of lags in the simulated ice-sheet/ice-shelf system and identify mechanisms by which the system responds to

atmospheric and oceanic forcings, we ran 32 sensitivity experiments (using the full grounding-line scheme) that isolated changes in air temperature, precipitation and sea surface temperature (ΔT_{air} , ΔP_{eff} and ΔSST) and simplified the time-varying RCP forcings (Fig. 3a, b). Each simplified forcing experiment comprised a perturbation-free interval from 0–2000 CE and a linear increase in T_{air} , P_{eff} or SST from 2000 CE to either 2100 CE or 2300 CE, after which the forcing was maintained unchanged for the remainder of the run (up to 5000 CE). This simplification was designed to clearly identify modelled ice-sheet response to each applied forcing.

Figure 3 illustrates the response of the modelled ice sheet when forced with each of the three environmental forcings (ΔT_{air} , ΔP_{eff} and ΔSST) in isolation, and when all three are combined, under the warmest of the simplified RCP scenarios (RCP 8.5) and using the full grounding-line scheme. Figure 3c and d shows that the timing of

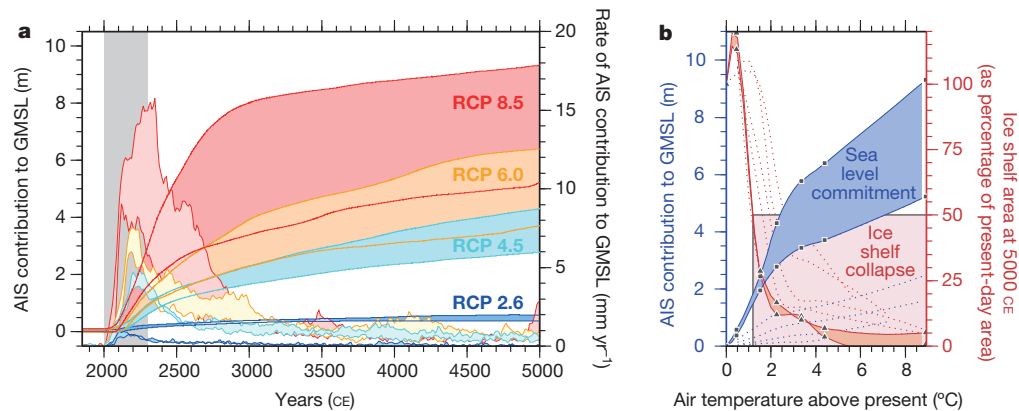


Figure 2 | Antarctic Ice Sheet (AIS) contribution to GMSL. **a**, Predicted sea-level contribution from the Antarctic ice sheet for 'high' and 'low' simulations (coloured lines) under each of the four RCP scenarios (darker shading), based on coeval climatic and oceanic perturbations. The forced response (grey shading) represents 20% to 36% of the committed response by 5000 CE. Lighter shading between coloured lines shows rates of sea-level-equivalent ice loss for each scenario. **b**, Long-term sea-level commitment as a function of atmospheric warming (blue shading with squares). Intermediate response

curves for the 'low' simulations are shown as dotted lines. Red shading with triangles shows the relationship between ice-shelf area and atmospheric warming for the near-equilibrium response and for intermediate stages (dotted lines). All curves in **b** are based on data from the four RCP scenario simulations, further constrained by two additional experiments whose maximum air temperature forcings are 1.5 °C and 3.35 °C. Pink shading defines the temperature range within which an ice shelf extent that is less than 50% of the present extent is simulated.

ice-sheet response to a warming atmosphere is similar for both 2100 and 2300 scenarios with rate-of-change maxima at 2244–2316 CE, several centuries from the start of the perturbation. In contrast, peak ice-volume responses to precipitation changes manifest only 11–41 years after the end of each forcing period. A warming ocean results in a 26-year peak-response lag when forced to 2100 CE, but there is no lag when forced to 2300 CE. In summary, precipitation changes produce rapid sea-level-equivalent ice-volume changes because of their direct and immediate addition of mass to the ice sheet. Warming of the atmosphere leads to an immediate loss of mass through surface melt in some areas (especially West Antarctica), but the effect is small and generally a warming atmosphere produces a lagged volumetric response that requires 200–300 years from the initial forcing to reach a maximum. Thermal changes in the ocean bring about rapid ice-sheet responses that, using the full grounding-line scheme, are far greater in magnitude than those arising from atmospheric forcings alone, and

exhibit little to no lag. Despite these differences, all forcings lead to committed responses that manifest as rates of change that are higher than those of the initial conditions for thousands of years after the forcing period (Extended Data Fig. 5).

The mechanisms responsible for the differing responses seen in Fig. 3 can be inferred from time-series of glaciological changes (Fig. 4). Rising air temperatures increase the total volume of temperate ice (ice that is close to melting point; brown line in Fig. 4a), which, because it is softer and more easily deformed, leads to an increase in the non-sliding, or 'creep', component of domain-averaged grounded-ice velocity (Fig. 4c) (the domain is the model grid). Together with an accompanying increase in sliding velocity, atmospheric thermal forcing leads to gradual thinning of parts of the ice sheet and a consequent reduction in the domain-averaged gravitational driving stress (Fig. 4b). Although the warming perturbation in each simulation is applied only from 2000–2300 CE, softening of the ice and non-sliding

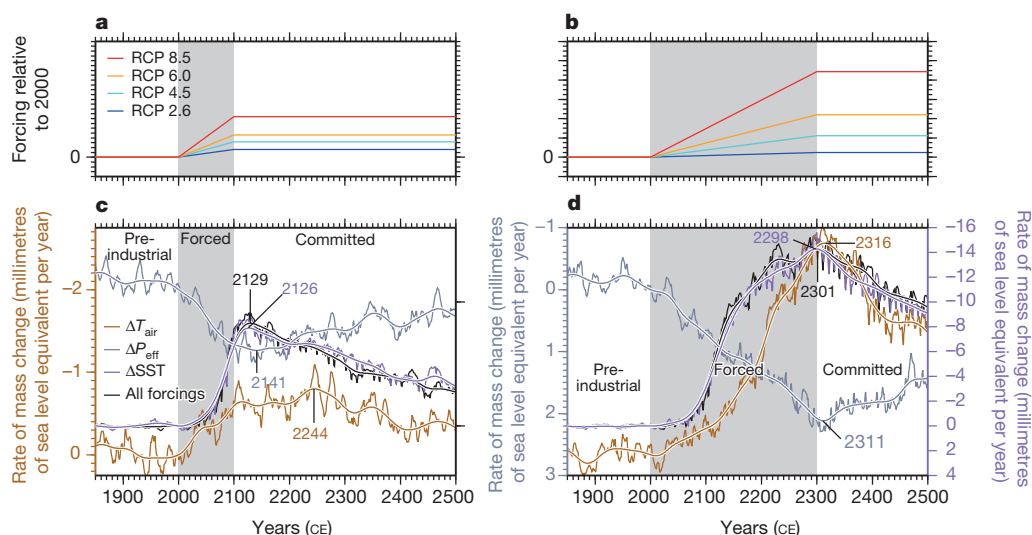


Figure 3 | Response of modelled Antarctic ice sheet to simplified environmental forcings. **a**, **b**, The relative magnitudes of air temperature ΔT_{air} , effective precipitation ΔP_{eff} and sea surface temperature ΔSST forcings for simplified RCP scenarios based on likely Antarctic values at 2100 CE (**a**) and 2300 CE (**b**). **c**, **d**, Trends in sea-level-equivalent ice mass rate-of-change that arise using the full grounding-line scheme and when forced with changes in ΔT_{air} , ΔP_{eff} , ΔSST , or the combination of these forcings, based on simplified

RCP 8.5 scenarios for 2100 (**c**) and 2300 (**d**). Note the different y-axis scales. Coloured numbers in **c** and **d** identify peak rates of change in the 100-year Gaussian-smoothed response curves. Data are shown relative to zero at 2000 CE. Grey shading shows periods of applied forcing. Absolute changes in ice-sheet volume and area (rather than rates of change) to 5000 CE are shown in Extended Data Fig. 4.

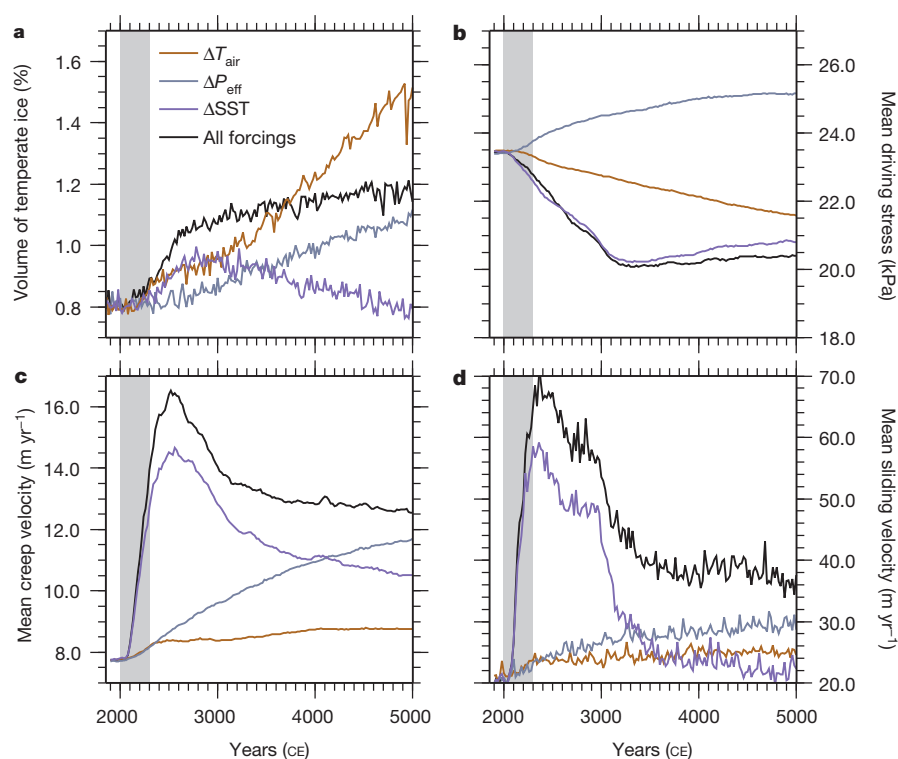


Figure 4 | Glaciological changes taking place under single-parameter environmental perturbations. **a**, Domain-wide changes in the volume of temperate ice. **b**, Mean driving stress of grounded ice. **c**, **d**, Averages of the non-sliding (**c**) and sliding (**d**) components of grounded ice velocity. All curves are from the simplified RCP 8.5, 2300 CE, 'high' simulations (Fig. 3b). Grey shading denotes period of applied forcing.

ice velocity continue to increase throughout the experiment, perhaps as a consequence of internal feedbacks such as strain heating, which would encourage faster ice flow and advection of 'warmer' ice deeper into the ice sheet, as proposed previously²⁸. Increased precipitation, even when applied in the absence of atmospheric warming, leads to an increase in the volume of temperate ice (grey line in Fig. 4a), probably because highest precipitation occurs in coastal areas under relatively warm conditions. Increased accumulation results in a greater domain-averaged ice thickness, which in turn increases the mean driving stress (Fig. 4b). Non-sliding ice velocities increase markedly throughout the simulation (Fig. 4c), probably as a consequence of increased ice thickness and increased driving stress²⁹.

By reducing the volume of floating ice shelves through basal melting, increased ocean temperatures lead to an overall reduction in the proportion of temperate ice in the domain (purple line in Fig. 4a). Thinning of grounded ice in response to loss of buttressing ice shelves means that the domain-averaged driving stress is reduced (Fig. 4b), but both sliding and non-sliding components of ice velocity increase abruptly during the forcing period (Fig. 4c, d). The acceleration of sliding velocities probably reflects the reduction in back-pressure exerted by the rapidly diminishing ice shelves⁷, whereas deformation rates probably increase as a consequence of internal strain heating and frictional heating at the bed due to faster flow²⁸. Flow acceleration then declines from about 3000 CE, perhaps because grounding lines occupy stable slopes once the major deep basins have been vacated, illustrating that the dynamic response to a loss of buttressing is transient and may be topographically controlled.

In summary, oceanic warming produces a much greater and more rapid ice-sheet response than atmospheric warming, but the rate of ice volume change forced by the warming of the ocean peaks and declines more quickly than the rate of change forced by atmospheric warming, which is more gradual and sustained for longer (Extended Data Fig. 5). Whereas ocean-forced perturbations produce the greatest ice-sheet contribution to sea level on centennial timescales (approximately five times that produced by increased air temperature alone when the full grounding-line scheme is used), contributions due to atmospheric warming become increasingly important over multi-millennial periods. The loss of buttressing ice shelves, the consequent increase in

sliding velocity of grounded ice and thinning at the ice-sheet margin, combined with the thermal softening effect of warming atmospheric temperatures that increases the creep-rate of grounded ice, all govern the long-term response of the ice sheet and lag times of this system. Ice-sheet response to external forcing is therefore mediated both by ice-stream dynamics and the timescale of adjustment of viscous flow ('creep') of grounded ice, which together produce a commitment to sea level rise that persists for multiple millennia.

Unless anthropogenic greenhouse-gas emissions are reduced to half of 1990 levels by 2050, global mean annual surface temperatures are likely to exceed 2 °C above pre-industrial values by 2100 (ref. 30). This aggressive mitigation scenario corresponds to RCP 2.6, which both our 'low' and 'high' simulations show is the only ocean-climate regime in which the long-term Antarctic contribution to GMSL does not exceed 1 m (Fig. 2a). Under all other RCP scenarios the future commitment to a rise in sea level from Antarctica is substantial. This commitment arises because the collapse of buttressing ice shelves leads to a dynamic ice-sheet response that greatly increases grounded ice discharge for hundreds to thousands of years into the future, even if greenhouse-gas emissions are reduced and temperatures stabilize (Supplementary Videos 1 and 2). The results presented here and elsewhere⁸ suggest that ice-shelf stability is vulnerable to a critical temperature threshold. In our experiments, we find that prolonged ocean warming of 0.5 °C above present, together with atmospheric warming of 2 °C, ultimately leads to the loss of 80% to 85% of all floating ice in Antarctica.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 13 April; accepted 3 September 2015.

- Meinshausen, M. *et al.* The RCP greenhouse gas concentrations and their extensions from 1765 to 2300. *Clim. Change* **109**, 213–241 (2011).
- Rogelj, J., Meinshausen, M. & Knutti, R. Global warming under old and new scenarios using IPCC climate sensitivity range estimates. *Nature Clim. Change* **2**, 248–253 (2012).
- Collins, M. *et al.* in *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* (eds Stocker, T. *et al.*) 1029–1136 (Cambridge Univ. Press, 2013).

4. Pritchard, H. D. *et al.* Antarctic ice-sheet loss driven by basal melting of ice shelves. *Nature* **484**, 502–505 (2012).
5. Wouters, B. *et al.* Dynamic thinning of glaciers on the Southern Antarctic Peninsula. *Science* **348**, 899–903 (2015).
6. Joughin, I. & Alley, R. B. Stability of the West Antarctic ice sheet in a warming world. *Nature Geosci.* **4**, 506–513 (2011).
7. Scambos, T. A., Bohlander, J. A., Shuman, C. A. & Skvarca, P. Glacier acceleration and thinning after ice shelf collapse in the Larsen B embayment, Antarctica. *Geophys. Res. Lett.* **31**, L18402 (2004).
8. Vaughan, D. G. & Doake, C. S. M. Recent atmospheric warming and retreat of ice shelves on the Antarctic Peninsula. *Nature* **379**, 328–331 (1996).
9. Liu, Y. *et al.* Ocean-driven thinning enhances iceberg calving and retreat of Antarctic ice shelves. *Proc. Natl Acad. Sci. USA* **112**, 3263–3268 (2015).
10. Paolo, F. S., Fricker, H. A. & Padman, L. Volume loss from Antarctic ice shelves is accelerating. *Science* **348**, 327–331 (2015).
11. Roemmich, D. *et al.* Unabated planetary warming and its ocean structure since 2006. *Nature Clim. Change* **5**, 240–245 (2015).
12. Schmidt, S., Heywood, K. J., Thompson, A. F. & Aoki, S. Multidecadal warming of Antarctic waters. *Science* **346**, 1227–1231 (2014).
13. Steig, E. J. *et al.* Warming of the Antarctic ice-sheet surface since the 1957 International Geophysical Year. *Nature* **457**, 459–462 (2009).
14. Harig, C. & Simons, F. J. Accelerated West Antarctic ice mass loss continues to outpace East Antarctic gains. *Earth Planet. Sci. Lett.* **415**, 134–141 (2015).
15. Hay, C. C., Morrow, E., Kopp, R. E. & Mitrovica, J. X. Probabilistic reanalysis of twentieth-century sea-level rise. *Nature* **517**, 481–484 (2015).
16. Joughin, I., Smith, B. E. & Medley, B. Marine ice sheet collapse potentially under way for the Thwaites Glacier basin, West Antarctica. *Science* **344**, 735–738 (2014).
17. Rignot, E., Mouginot, J., Morlighem, M., Seroussi, H. & Scheuchl, B. Widespread, rapid grounding line retreat of Pine Island, Thwaites, Smith, and Kohler glaciers, West Antarctica, from 1992 to 2011. *Geophys. Res. Lett.* **41**, 3502–3509 (2014).
18. Bindenschadler, R. *et al.* Ice-sheet model sensitivities to environmental forcing and their use in projecting future sea-level (The SeaRISE Project). *J. Glaciol.* **59**, 195–224 (2013).
19. Levermann, A. *et al.* Projecting Antarctic ice discharge using response functions from SeaRISE ice-sheet models. *Earth Syst. Dyn.* **5**, 271–293 (2014).
20. Levermann, A. *et al.* The multimillennial sea-level commitment of global warming. *Proc. Natl Acad. Sci. USA* **110**, 13745–13750 (2013).
21. Naish, T. *et al.* Obliquity-paced Pliocene West Antarctic ice sheet oscillations. *Nature* **458**, 322–328 (2009).
22. Pollard, D. & DeConto, R. M. Modelling West Antarctic ice sheet growth and collapse through the past five million years. *Nature* **458**, 329–332 (2009).
23. Schaeffer, M., Hare, W., Rahmstorf, S. & Vermeer, M. Long-term sea-level rise implied by 1.5 °C and 2 °C warming levels. *Nature Clim. Change* **2**, 867–870 (2012).
24. Bueler, E. & Brown, J. Shallow shelf approximation as a “sliding law” in a thermomechanically coupled ice sheet model. *J. Geophys. Res.* **114**, F03008 (2009).
25. Winkelmann, R. *et al.* The Potsdam Parallel Ice Sheet Model (PISM-PIK)—Part 1: Model description. *Cryosphere* **5**, 715–726 (2010).
26. Feldmann, J. & Levermann, A. Interaction of marine ice-sheet instabilities in two drainage basins: simple scaling of geometry and transition time. *Cryosphere* **9**, 631–645 (2015).
27. Feldmann, J., Albrecht, T., Khroulev, C., Pattyn, F. & Levermann, A. Resolution-dependent performance of grounding line motion in a shallow model compared to a full-Stokes model according to the MISIP3d intercomparison. *J. Glaciol.* **60**, 353–360 (2014).
28. Clarke, G. K., Nitsan, U. & Paterson, W. Strain heating and creep instability in glaciers and ice sheets. *Rev. Geophys.* **15**, 235–247 (1977).
29. Winkelmann, R., Levermann, A., Frieler, K. & Martin, M. Increased future ice discharge from Antarctica owing to higher snowfall. *Nature* **492**, 239–242 (2012).
30. Meinshausen, M. *et al.* Greenhouse-gas emission targets for limiting global warming to 2 °C. *Nature* **458**, 1158–1162 (2009).

Supplementary Information is available in the online version of the paper.

Acknowledgements We thank the CMIP community for making their data openly available, and J. Lenaerts for providing present-day surface mass balance data. We are also grateful to K. Buckley (Victoria University high-performance computing cluster), C. Khroulev, T. Albrecht and the Parallel Ice Sheet Model groups at the University of Alaska, Fairbanks, and the Potsdam Institute for Climate Impact Research. This work was funded by contract VUW1203 of the Royal Society of New Zealand’s Marsden Fund, with support from the Antarctic Research Centre, Victoria University of Wellington, ANDRILL, GNS Science (NZ Ministry of Business Innovation and Employment contract C05X1001), National Science Foundation grant ANT-1043712, and the Australian Research Council (ARC). J. Renwick and D. Zwart provided comments that improved the manuscript.

Author Contributions N.R.G. devised and carried out the ice-sheet modelling experiments and D.E.K. undertook climate model simulations to produce the present-day ocean temperature field. All authors contributed to the development of ideas and writing of the manuscript.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to N.R.G. (nicholas.golledge@vuw.ac.nz).

METHODS

The ice-sheet model. Continental-scale studies such as ours are currently unable to reproduce the level of detail that catchment-scale studies resolve^{31,32}, but they nonetheless offer useful insights that arise from their wider geographic coverage. For our suite of experiments we use the Parallel Ice Sheet Model (PISM) version 0.6.3, an open-source, three-dimensional, thermodynamic, coupled ice-sheet/ice-shelf model. The model combines equations of the shallow-ice and shallow-shelf approximations for grounded ice, and uses the shallow-shelf approximation for floating ice²⁴. Superposing the shallow-ice and shallow-shelf velocity solutions allows basal sliding to be simulated according to the ‘dragging shelf’ approach²⁴, and enables a consistent treatment of stress regime across the grounded-ice to floating-ice transition²⁵. Ice streams develop naturally as a consequence of plastic failure of saturated basal till³³, depending on the thermal regime and volume of water at the ice-sheet bed. The amount of water that saturates basal till varies spatially according to the enthalpy field, but is limited in our implementation so that no more than a 2-m thickness of saturated substrate accumulates. This sub-glacial hydrology model does not conserve water, in the sense that any additional water above the imposed maximum thickness is permanently lost. Till pore water pressure is defined as a maximum of 0.8 times the ice-overburden pressure (Extended Data Table 2). The effective pressure calculated from meltwater thickness and ice overburden, together with the spatially varying till friction angle, is then used to calculate the yield strength of basal substrate, which will vary in time and space as basal meltwater volumes change. The till friction angle is prescribed heuristically at the start of the run, and follows a simple elevation dependence such that deeper basins that are likely to hold deformable sediments are assigned lower values than higher areas, which are more likely to be rock. In our simulations, we choose values of 10° for areas below −1,000 m, and 30° for areas above 200 m. Intermediate values are linearly interpolated between these extrema. For computational tractability we employed a resolution of 10 km for our ‘realistic’ RCP simulations, and a resolution of 20 km for our ‘simplified’ RCP experiments, where the focus lies on investigating the differences between results of different environmental forcings, rather than simply the magnitudes or rates of ice retreat.

Migration of the grounding line is facilitated through a sub-grid scheme that calculates one-sided derivatives for surface slope above and below the grounding line, and interpolates basal shear stress in x , y according to the spatial gradient between adjacent grounded and floating cells²⁷. Additionally, we impose reduced basal traction in the first cell upstream of the grounding line, in which the basal resistance calculated from the computed pore water content is replaced by a basal drag value that assumes saturation of the basal substrate. In this way, basal shear stress gradients across the grounding line are reduced and a more dynamic margin is facilitated, in line with current theory^{34,35}. For our ‘high’ simulations we also use a sub-grid scheme that interpolates sub-shelf melt rates at the grounding line^{26,27}.

Climate data (air temperature, precipitation) are used as inputs to a positive degree-day scheme that calculates surface mass balance. We use degree-day factors of 3 mm K^{−1} per day and 8 mm K^{−1} per day for snow and ice respectively. Random white noise is included in the calculation to mimic natural variability. An elevation dependence of −0.008 K m^{−1} is applied to air temperatures in areas where the surface elevation changes. Basal melting of floating ice is calculated thermodynamically using a three-component temperature–salinity scheme³⁶. This scheme yields the highest melt rates adjacent to ice-sheet grounding lines and lower melt rates, or accretion, closer to the central sectors of modelled ice shelves. Calving of floating ice is calculated from horizontal strain rates^{37,38} with an additional heuristically based minimum thickness condition applied. PISM incorporates a sophisticated bed deformation model in which changes in ice load effect both elastic and viscous responses in the underlying bedrock³⁹. Time-dependent variations in relative sea level therefore arise in our simulations from isostatic adjustment according to the flexural rigidity of the crust and the viscosity-dependent lateral displacement of the mantle. Self-gravitational sea-level effects that arise from changes in ice-sheet mass are not included, however, and consequently our model may not capture the potentially stabilizing feedback of a local lowering of sea level⁴⁰. This effect, which becomes more important over multi-centennial timescales, is typically smaller than isostatic effects, which are included in our simulations, so the influence on our short-term predictions should be minimal. Because we also do not account for self-gravitational rises in local sea level, which may occur in response to loss of ice from distal sectors of the ice sheet, we may in fact underestimate margin retreat in some areas.

Climate and ocean inputs.

Present-day conditions. To define the initial conditions for our experimental ensemble, we employ spatially distributed data sets of annual mean air temperatures, annual precipitation totals, and circum-Antarctic sea surface temperatures from a selection of recent observational and model-based compilations.

The most recent and widely used map of Antarctic surface mass balance derives from regional atmospheric climate modelling using RACMO2.1/ANT45. These

data span the period 1979–2010, are high-resolution (27 km), and incorporate snow-drift physics that greatly improve the fit of simulated mass balance to around 750 *in situ* measurements (compared to models that ignore snow drifting). By capturing the erosional as well as the sublimation effects of wind, the model⁴¹ is able to compute both surface accumulation and the location and extent of ablation areas. Annual mass balance across the continent is calculated to be $2,418 \pm 181$ Gt per year with only minor interannual variability but a pronounced, winter-dominated, seasonal cycle.

Antarctic surface air temperatures derived from Advanced Very High Resolution Radiometer (AVHRR) infrared data⁴² and updated for the period 1982–2004 (ref. 43) provide high-resolution (6.25 km) continent-wide coverage. Since temperatures are more easily and reliably mapped than surface accumulation, these data^{42,43} are commonly used in ice-sheet modelling.

Oceanic temperatures are harder to measure, and for our experiments we use outputs from an established regional-scale climate model. This model (RegCM3) is coupled with the GENESIS version 3.0 Global Climate Model and is implemented at 80-km resolution^{44,45}. The model includes parameterizations of surface, boundary layer and most processes that account for the physical exchanges between the land surface boundary layer and free atmosphere. It includes detailed representations of snow and near-surface land ice and has been validated against modern observed polar climates and ice-sheet mass balances⁴⁴. Oceanic outputs from this model capture the continental-scale pattern of Antarctic sea surface temperatures but probably underestimate temperatures in areas where recent subsurface warming has been most rapid, such as the inner Amundsen Sea. The underestimation of modern ocean temperatures in this area may be the primary reason that our simulations do not appear to show rapid grounding-line retreat into the Pine Island and Thwaites Glacier subglacial basins, as has been predicted by some models¹⁶, but we also note that the muted and lagged response we see in the Amundsen Sea sector is consistent with the findings of another study⁴⁶.

Since our primary interest lies in continental-scale interpretations, and since our methodology is based on the analysis of deviations from initial conditions (that is, results are bias-corrected), we argue that if any local inaccuracies in our input fields exist they should not substantially affect our conclusions at the continental scale; however, we acknowledge this as an area where future refinements could be made.

Future climate trajectories. To perturb our initial conditions along the trajectories forecast for coming centuries, we use Climate Model Intercomparison Project phase 5 (CMIP5) data for the four RCPs of the Fifth Assessment Report of the Intergovernmental Panel on Climate Change³. We extracted Antarctic-specific (60°–90° S) zonally averaged timeseries of air temperature (ΔT_{air}), effective precipitation (ΔP_{eff}) and sea surface temperature (ΔSST) from the ensemble mean data set for the period 1890–2100 CE (Extended Data Fig. 1 and Extended Data Table 1). Because we use the ensemble mean, rather than the full suite of climate model projections presented in the Fifth Assessment Report of the Intergovernmental Panel on Climate Change³, our simulations do not capture the full range of possible responses implied by the spread of climate model predictions. However, our aim is to present representative outputs that allow first-order differences between the main scenarios to be identified.

Raw data were uniformly adjusted so that perturbations are relative to the present, to ensure consistency with our initial-condition data sets (above). This adjustment effectively results in air temperatures at 1890 CE approximately 0.6 °C below those of the present³. For each of the RCPs, we used the extracted Antarctic-specific air temperature, precipitation and ocean temperature anomalies to update the environmental boundary conditions of our ice-sheet model incrementally throughout the simulations. To extend the CMIP5 RCP values beyond 2100 CE we used long-term trajectories based on the Fifth Assessment Report of the Intergovernmental Panel on Climate Change^{3,47}. In these Extended Concentration Pathway scenarios, all emissions trajectories result in maximum atmospheric temperature perturbations by 2300 CE, and remain unchanged for the remainder of the simulation period (to 5000 CE)^{1,3,47} (Extended Data Fig. 2). In all of our experiments, the magnitudes of precipitation and ocean temperature changes are scaled to air temperature changes according to their respective variance in the zonally averaged CMIP5 data. These indicate that a 1 °C increase in air temperature accounts for a 5.3% increase in precipitation, in agreement with long-term records⁴⁸, whereas ocean temperature changes are approximately one-quarter of that seen in the atmosphere (Extended Data Fig. 1).

In chapter 12 of the Fifth Assessment Report of the Intergovernmental Panel on Climate Change³, long-term (beyond 2300 CE) climate change projections are considered in terms of two scenarios: first, a situation in which atmospheric CO₂ remains constant at 2300 CE levels, and second, where atmospheric CO₂ declines steadily from 2300 CE to 3000 CE. This second scenario, however, results in only a modest reduction in surface air temperatures, and is based on reduction of emissions to zero by 2300 CE. The likelihood of such aggressive mitigation is not

known, but the available data appear to “...support the conclusion that temperatures would decrease only very slowly (if at all), even from strong reductions or complete elimination of CO₂ emissions...” (page 1104 of ref. 3). On this basis, we propagate our environmental forcings beyond 2300 CE at their 2300 CE perturbed level. In most cases this is their maximum, although in RCP 2.6 maxima are reached by 2100 CE, after which they decline.

Code availability. The Parallel Ice Sheet Model is freely available as open-source code from the PISM github repository (<https://github.com/pism/pism>). RegCM3 is available from <https://users.ictp.it/RegCNET/model.html>. CMIP5 data were downloaded from <http://climexp.knmi.nl/>. Bedrock topography and ice thickness data are from the BEDMAP2 compilation, available at <http://www.antarctica.a-c.uk/basresearch/ourresearch/az/bedmap2/>. Information on surface mass balance data is available at <http://www.projects.science.uu.nl/iceclimate/models/antarctica.php#racmo21>. Air temperature and geothermal heat flux inputs were taken from the ALBMAP version 1 compilation⁴³ and can be downloaded from <http://doi.pangaea.de/10.1594/PANGAEA.734145>.

Experimental methods. To establish differences in ice-sheet configuration arising under perturbed climate regimes, our initial requirement is for an accurate simulation of the present-day Antarctic ice sheet in which the grounded and floating portions are reproduced in a manner that is as close to observations as possible. Our principal guiding constraints are (1) the volume of grounded ice, (2) grounding-line positions, and (3) the pattern of surface velocities. Building on previous work^{49,50} we implement the following three-stage spin-up procedure in order to produce a thermally equilibrated and dynamically stable simulation of the present-day ice sheet, from which subsequent experiments can be restarted:

(1) Initial 20-year smoothing run in which only the shallow-ice approximation is used to calculate ice flow. This allows the initial ice-sheet surface to relax slightly, removing any anomalies introduced in the initial data collation phase. During this brief run, the calving line is held fixed.

(2) Intermediate 150,000-year run in which the ice geometry is held fixed, but the enthalpy field is allowed to evolve. This allows three-dimensional ice temperatures to evolve to equilibrium under the imposed initial climate conditions.

(3) Final 25,000-year run in which full model physics are employed, including both shallow-ice and shallow-shelf approximations for velocity calculations²⁴, viscoelastic bed deformation³⁹, grounding-line migration²⁷ and calving^{37,38}.

During the 25,000-year final part of this spin-up procedure, all model boundaries (calving line, grounding line, upper and lower surfaces) are free to evolve. Consequently, achieving an ice-sheet geometry at the end of the spin-up that is close to the observed present-day configuration requires careful parameterization. Iterative experimentation focused on manually adjusting flow enhancement factors, basal hydrology and basal traction parameters, and calving coefficients. For each experiment ice thickness and surface velocities were compared to observation-based data to gauge the degree of fit. Because we do not use inverse or iterative methods⁵¹ our fit to present-day constraints has outliers (Extended Data Fig. 3). However, the positions of grounding lines and calving lines are well captured, and the sea-level-equivalent ice volume is within 4% of empirically calculated values⁵². By allowing our tuning experiments to evolve over such long timescales we ensure both that our results are not influenced by transient behaviours and also that our parameterization is robust for multi-millennial integrations. This is essential, since to identify climate-forced perturbations in our study, model drift from incorrect parameterization needs to be minimized.

From the optimal parameterization spin-up run we obtain a thermally equilibrated and dynamically stable present-day ice-sheet simulation that closely resembles the modern ice sheet (Extended Data Fig. 3). This simulation is the starting point for all of the RCP-based experiments. We then run a two-component experimental ensemble in which each of the environmental parameters (air temperature T_{air} , precipitation P_{eff} , ocean temperature SST) evolves according to the rates and magnitudes indicated in the extended CMIP5 ensemble mean trends as described above. One component of our experimentation involves 10-km-resolution runs using the full RCP forcing scenarios. The other component uses 20-km-resolution runs to explore glaciological changes that occur in response to theoretical environmental perturbations. Each of the experiments runs for 5,000 years, starting at year zero with climate and ocean forcings applied in calendar years through to 5000 CE. To minimize the possibility of any transient effects arising during the start-up procedure, environmental boundary conditions are held constant for the first 2,000 years, but all model boundaries (calving line, grounding line, upper and lower surfaces) are free to evolve. Any transient adjustments of the ice sheet occur during this period. From 2000 CE, timeseries environmental forcings (described above) are imposed and the ice-sheet/ice-shelf system evolves accordingly. In addition to the forcing experiments we run an additional control experiment in which the same 5,000-year run is made, starting from identical initial conditions but with no environmental perturbations applied. This is essential for the identification and quantification of any model drift during

the period of interest of the simulation (2000–5000 CE). Since this drift must arise from processes other than environmental forcing we use this control experiment to bias-correct outputs from all other simulations. Even though all model boundaries are allowed to evolve freely, model drift in our control experiment is near zero during the period of interest for both ice-volume and ice-sheet extent (Extended Data Fig. 4), so we do not consider that the bias-correction affects our results. In fact, we consider the low model drift to be additional confirmation that our parameterization is robust. On this basis we are confident that the response seen in the RCP simulations arises solely as a consequence of the combined influence of the simultaneously applied forcings (Fig. 1).

Grounding-line sensitivity and resolution dependence. A key part of our simulations is the migration of ice-sheet grounding lines in response to applied environmental forcings. In theory, to obtain the most accurate predictions of grounding-line locations a numerical scheme is required that is able to solve the full set of Stokes equations for ice flow on a high-resolution (hundreds to thousands of metres) spatial grid⁵³. For continental-scale simulations, however, such an approach presents large challenges, and therefore various alternative schemes exist^{22,27}. These alternatives attempt to capture the critical aspects of grounding-line behaviour while still remaining tractable over large domains (thousands of square kilometres) and long (multi-millennial) model integrations. Of primary concern is whether the grid size over which the scheme is implemented actually influences the model results. Here we use two new grounding-line components that together allow dynamic grounding-line behaviour even with our relatively coarse (10 km) grid, but which may also be regarded as sources of uncertainty in our future ice-volume projections.

The first component is a sub-grid grounding-line scheme²⁷ that interpolates basal shear stress (and, optionally, basal melt fields) at the grounding line in order to facilitate the smoother migration of the grounding line. This approach attempts to capture the effects of lateral contact between ice and ocean as well as the effects of water intrusion beneath the ice sheet. Furthermore, it is logical that a retreating grounding line must be provided with an (oceanic) basal melt value at the sub-grid cell into which it migrates, rather than a basal melt value typical of the grounded ice that characterizes the uninterpolated cell. Whether or not the grounding line would migrate without this interpolation, however, remains a source of uncertainty. Thus in our simulations we present both ‘high’- and ‘low’-scenario estimates of sea-level-equivalent ice loss based on turning on or off the sub-grid melt interpolation. In other studies that use this scheme⁴⁶, some evidence of resolution dependency is apparent, but whether this is a consequence of the sub-grid scheme or of some other aspect of the model implementation is not clear.

The second component of the grounding-line scheme we adopt is an heuristic scheme that reduces basal traction in the first cell upstream of the grounding line. Smoothing the basal friction gradient across the grounding line should reduce the grid-dependency problems seen in artificial scenarios such as MISMP3d^{27,54} that invoke large changes in basal traction between grounded and floating domains. However, smoothing the basal friction gradient at the grounding line may make this critical boundary more sensitive to environmental perturbations than in models where basal traction changes abruptly, and thus this also perhaps introduces an element of uncertainty in our ice-volume projections.

Grounding lines are poorly observed in nature due to their inaccessibility, yet ice-shelf basal melt rates typically reach their maxima close to the grounding line, and may be as high as 100 m per year (ref. 55). These rates are, however, highly variable and thus while small environmental perturbations could result in relatively large ice-sheet responses in some areas, the consequences in other sectors could be more muted. Acknowledging these uncertainties, we ran a series of experiments at 10-km resolution in an attempt to quantify (at the continental scale and at the multi-millennial scale) the effects of removing individual aspects of the grounding-line scheme. Extended Data Fig. 6a illustrates sea-level-equivalent volume changes for each of these experiments, in which all other parameters were held constant. Using the sub-grid scheme without the basal traction scheme (green line in Extended Data Fig. 6a) results in a relatively minor change to ice volume compared to the RCP 8.5 simulation that uses the full grounding-line scheme (blue line), whereas using the basal traction scheme without the sub-grid scheme (purple line) results in a much less dynamic ice margin that responds less to applied environmental forcings. Using the basal traction scheme together with the sub-grid scheme in which basal melt is not interpolated produces an intermediate result (grey line). Extended Data Fig. 6e–g illustrates the effect of this change in terms of grounding-line migration through time. If neither the sub-grid nor basal traction schemes are used, the simulated ice sheet not only increases in volume above present-day values, but also becomes less responsive to the applied forcings (orange line).

This sub-grid grounding-line scheme is clearly important, and has been shown to greatly improve the accuracy of grounding-line migrations compared to earlier PISM releases⁵⁴, but its implementation may still result in a certain amount of

resolution dependence under certain circumstances^{27,46}. To assess the uncertainties that may arise from grid resolution alone, we implemented a series of experiments in which we employed identical parameterizations (including sub-grid melt interpolation) except for the horizontal grid spacing used to solve prognostic equations. Extended Data Fig. 6b and c illustrates the convergence of results that occurs when the grid is progressively refined from 80 km to 5 km. The finest of these runs (5 km) was conducted only to 3100 CE, owing to the large computational overhead. In terms of final sea-level-equivalent ice volumes, there is a greater difference between the 40-km and 80-km runs than there is between the 20-km and 40-km runs, and there is no difference at all between the final state of the 10-km and 20-km simulations. At even finer grids, the overall trajectories of ice-volume and grounded-ice area at 10 km resolution agree closely with those of the first three millennia of the 5-km-resolution simulation. The spatial variability that arises at these different resolutions is illustrated in Extended Data Fig. 6d for the 2500 CE time slice, identifying that the greatest resolution-dependent uncertainty appears to occur in the Siple Coast region of West Antarctica, but even in this area there is little difference between the 10-km and 5-km results after 500 years of environmental forcing.

Although the 10-km and 20-km simulations converge on identical final ice volumes, the coarser of these two experiments actually simulates slightly greater ice loss, because it deviates from the 'control' ice volume of our simulated present-day ice sheet by approximately ± 0.5 m sea-level equivalent. To see whether this offset could be reduced, we ran 20 new experiments that explored a range of stress balance parameterizations. We found that by making minor adjustments to both the shallow-ice and shallow-shelf flow enhancement factors (Extended Data Table 2), fully dynamic simulations of the present-day ice sheet could be produced (at either resolution) that are freely evolving over multi-millennial timescales but exhibit minimal model drift and fit equally well to present-day geometric and dynamic constraints at the continental scale. Although these parameterizations were tuned only to reproduce the geometry and dynamics of the present-day ice sheet, we found that similar patterns and rates of grounding-line migration were simulated under RCP forcing scenarios (Extended Data Fig. 6h–j).

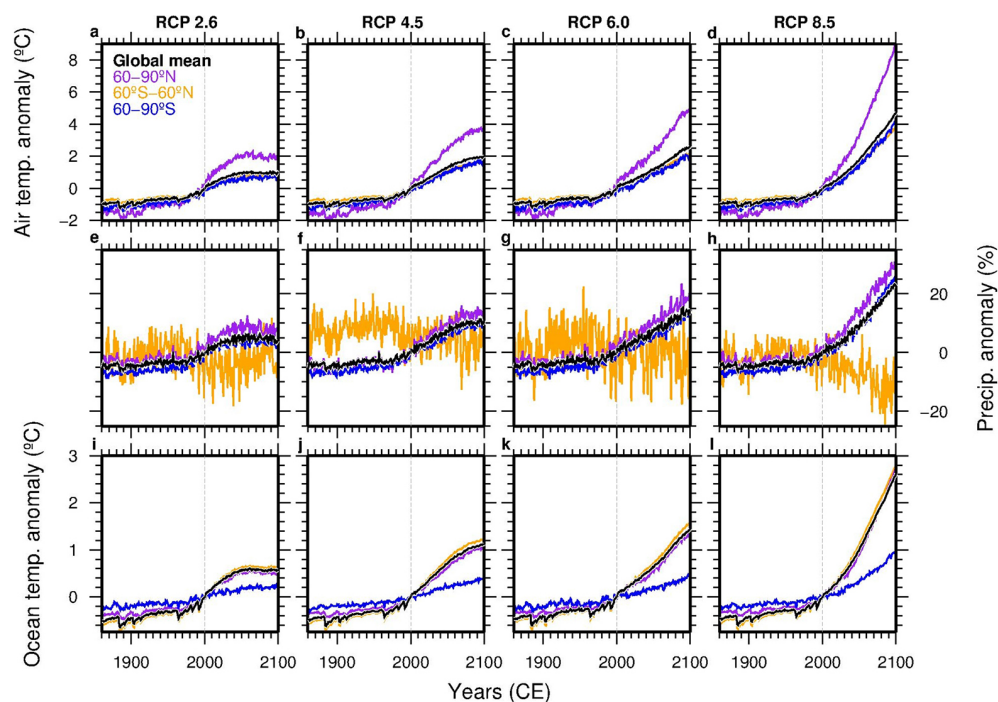
Together, the suite of sensitivity experiments described above illustrate that grounding-line treatment can strongly influence predictions of both the rate of ice-sheet retreat, and the total volume lost for any given forcing. Although the sub-grid melt interpolation scheme still requires physical verification, and tends to favour grounding-line retreat, values from our 'high' simulations may still underestimate ice-sheet responses that arise from retreat mechanisms or systemic feedbacks not included in our model. Likewise, the more conventional implementation of sub-grid grounding-line scheme that neglects interpolated basal melt results in more conservative sea-level contributions in our 'low' simulations, but still lower values could be possible if unaccounted-for stabilizing mechanisms come into play. Given that geological reconstructions tend to indicate a relatively high ice-sheet sensitivity to warmer climates of the past^{3,56}, we find no compelling reason to favour the conventional grounding-line implementation over the one that includes interpolated basal melt. Regardless of the absolute magnitudes, however, our simulations consistently identify a three-stage response to environmental warming: an initial loss or reduction in ice-shelf area, a dynamic response of grounded ice, and a long-term commitment to sea-level rise that continues for centuries to millennia beyond the initial forcing period.

Polar amplification and long-term equilibrium response. Although there exists some uncertainty as to the form of the RCP projections, particularly beyond 2300 CE, in our experiments we are concerned with the range of likely responses to a plausible range of environmental forcings, and as such we are confident that our ensemble approach captures the likely breadth of Antarctic ice sheet response. Palaeoclimate data and models indicate that, over millennial timescales, polar amplification processes can lead to surface air temperatures, and ocean temperatures, that may be as much as double the magnitude of the global mean perturbation^{56–58}. If this is the case, then our RCP-based experiments represent a conservative lower estimate of the likely ice-sheet response. Furthermore, because we do not use a coupled climate model in our simulations, we cannot fully capture elevation-dependent mass-balance feedbacks that might take place as the ice-sheet geometry evolves. Although we use an elevation-dependent lapse rate to modify air temperatures through the run, it is also likely that, where ice shelves and grounded ice in marine basins are greatly reduced in extent, circulation changes will arise that may bring warmer and wetter air masses further into East Antarctica. To assess the likely implications of even greater warming of the southern high latitudes, we therefore ran duplicate experiments of the RCP 8.5 scenario, with and without sub-grid melt interpolation, and allowed air and ocean temperature anomalies to double by 2300 CE, compared to the original RCP 8.5 experiment. This yields surface air and ocean temperature anomaly maxima of 16°C and 4°C respectively, by 2300 CE and beyond. The consequence of this extreme

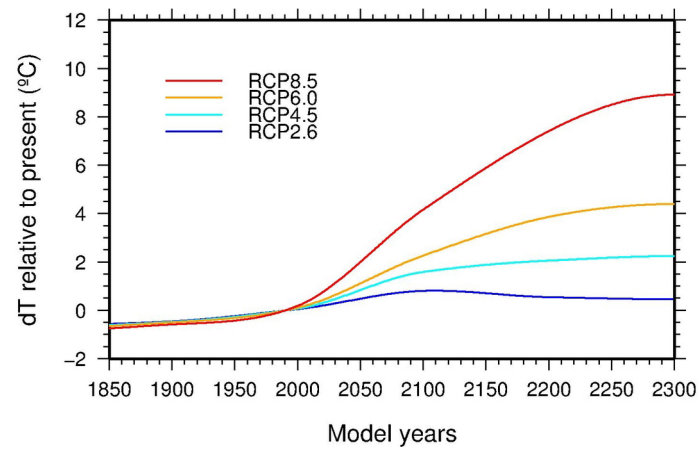
scenario compared to the original RCP 8.5 scenario is illustrated in Extended Data Figs 7 and 8.

We note that although the ice-sheet geometry is visually similar between RCP 8.5 and RCP 8.5A for the 'high' scenario (blue outline), the sea-level contribution is 2.1 m greater by 5000 CE under a polar-amplified climate, and the ongoing rate of ice mass loss is 80% higher (Extended Data Fig. 7a, b). In the 'low' experiments the difference in sea-level contribution is much greater, with the loss of the WAIS and a considerably higher sea-level contribution arising with amplified warming (8.57 m compared to 5.2 m). With polar amplification, grounding-line retreat into Wilkes Basin occurs under both 'high' and 'low' scenarios (Extended Data Fig. 7b, c). Extending the simulation to 50,000 years (Extended Data Fig. 7c, d) the full extent of the warming-based commitment under the 'high' scenario is far greater than that reached by 5000 CE (15.73 m sea-level equivalent compared to 11.42 m). Under the 'low' scenario the difference is much less (9.68 m compared to 8.57 m).

31. Favier, L. *et al.* Retreat of Pine Island Glacier controlled by marine ice-sheet instability. *Nature Clim. Change* **4**, 117–121 (2014).
32. Cornford, S. *et al.* Century-scale simulations of the response of the West Antarctic Ice Sheet to a warming climate. *Cryosphere* **9**, 1579–1600 (2015).
33. Schoof, C. A variational approach to ice stream flow. *J. Fluid Mech.* **556**, 227–251 (2006).
34. Leguy, G. R., Asay-Davis, X. S. & Lipscomb, W. H. Parameterization of basal friction near grounding lines in a one-dimensional ice sheet model. *Cryosphere* **8**, 1239–1259 (2014).
35. Tsai, V. C., Stewart, A. L. & Thompson, A. F. Marine ice-sheet profiles and stability under Coulomb basal conditions. *J. Glaciol.* **61**, 205–215 (2015).
36. Holland, D. M. & Jenkins, A. Modeling thermodynamic ice-ocean interactions at the base of an ice shelf. *J. Phys. Oceanogr.* **29**, 1787–1800 (1999).
37. Albrecht, T. & Levermann, A. Fracture field for large-scale ice dynamics. *J. Glaciol.* **58**, 165–176 (2012).
38. Levermann, A. *et al.* Kinematic first-order calving law implies potential for abrupt ice-shelf retreat. *Cryosphere* **6**, 273–286 (2012).
39. Bueler, E. D., Lingle, C. S. & Brown, J. Fast computation of a viscoelastic deformable Earth model for ice-sheet simulations. *Ann. Glaciol.* **46**, 97–105 (2007).
40. Gomez, N., Pollard, D., Mitrovica, J. X., Huybers, P. & Clark, P. U. Evolution of a coupled marine ice sheet–sea level model. *J. Geophys. Res.* **117**, F01013 (2012).
41. Lenaerts, J., van den Broeke, M., van de Berg, W., van Meijgaard, E. & Munneke, P. A new, high-resolution surface mass balance map of Antarctica (1979–2010) based on regional atmospheric climate modeling. *Geophys. Res. Lett.* **39**, L04501 (2012).
42. Comiso, J. Variability and trends in Antarctic surface temperatures from in situ and satellite infra-red measurements. *J. Clim.* **13**, 1674–1696 (2000).
43. Le Brocq, A., Payne, A. & Veli, A. An improved Antarctic dataset for high resolution numerical ice sheet models (ALBMAP v1). *Earth Syst. Sci. Data* **2**, 247–260 (2010).
44. Thompson, S. L. & Pollard, D. Greenland and Antarctic mass balances for present and doubled atmospheric CO₂ from the genesis version-2 global climate model. *J. Clim.* **10**, 871–900 (1997).
45. Pal, J. S. *et al.* Regional climate modeling for the developing world: the ICTP RegCM3 and RegCM3. *Bull. Am. Meteorol. Soc.* **88**, 1395–1409 (2007).
46. Martin, M. A., Levermann, A. & Winkelmann, R. Comparing ice discharge through West Antarctic Gateways: Weddell vs. Amundsen Sea warming. *Cryosphere Discuss.* **9**, 1705–1733 (2015).
47. Zickfeld, K. *et al.* Long-term climate change commitment and reversibility: an EMIC intercomparison. *J. Clim.* **26**, 5782–5809 (2013).
48. Frieler, K. *et al.* Consistent evidence of increasing Antarctic accumulation with warming. *Nature Clim. Change* **5**, 348–352 (2015).
49. Gollidge, N. R., Fogwill, C. J., Mackintosh, A. N. & Buckley, K. M. Dynamics of the Last Glacial Maximum Antarctic ice-sheet and its response to ocean forcing. *Proc. Natl Acad. Sci. USA* **109**, 16052–16056 (2012).
50. Gollidge, N. *et al.* Antarctic contribution to meltwater pulse 1A from reduced Southern Ocean overturning. *Nature Commun.* **5**, 1–10 (2014).
51. Pollard, D. & DeConto, R. M. A simple inverse method for the distribution of basal sliding coefficients under ice sheets, applied to Antarctica. *Cryosphere* **6**, 953–971 (2012).
52. Fretwell, P. *et al.* Bedmap2: improved ice bed, surface and thickness datasets for Antarctica. *Cryosphere* **7**, 375–393 (2013).
53. Durand, G., Gagliardini, O., Zwinger, T., Le Meur, E. & Hindmarsh, R. C. Full Stokes modeling of marine ice sheets: influence of the grid size. *Ann. Glaciol.* **50**, 109–114 (2009).
54. Pattyn, F. *et al.* Grounding-line migration in plan-view marine ice-sheet models: results of the ice2sea MISMIP3d intercomparison. *J. Glaciol.* **59**, 410–422 (2013).
55. Dutrieux, P. *et al.* Pine Island glacier ice shelf melt distributed at kilometre scales. *Cryosphere* **7**, 1543–1555 (2013).
56. Naish, T. & Zwartz, D. Palaeoclimate: looking back to the future. *Nature Clim. Change* **2**, 317–318 (2012).
57. Graversen, R. G. & Wang, M. Polar amplification in a coupled climate model with locked albedo. *Clim. Dyn.* **33**, 629–643 (2009).
58. Masson-Delmotte, V. *et al.* in *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* (eds Stocker, T. *et al.*) 383–464 (2013).
59. Rignot, E., Mouginot, J. & Scheuchl, B. Ice flow of the Antarctic Ice Sheet. *Science* **333**, 1427–1430 (2011).



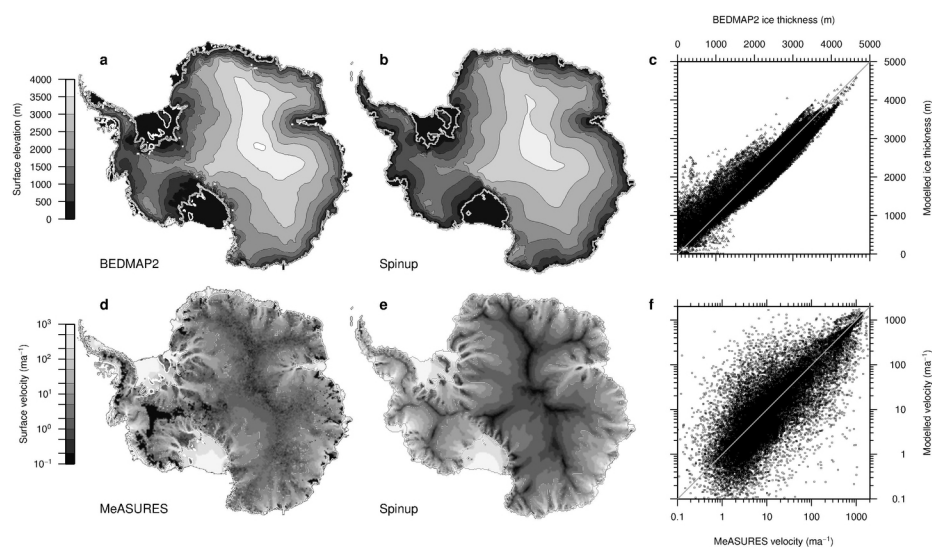
Extended Data Figure 1 | CMIP5 multi-model ensemble mean anomaly timeseries data. Air temperature (a–d), precipitation (e–h) and ocean temperature (i–l) changes for four RCP scenarios expressed as perturbations from present, both for hemispheric sectors and for the global mean.



Extended Data Figure 2 | Long-term RCP temperature scenarios.

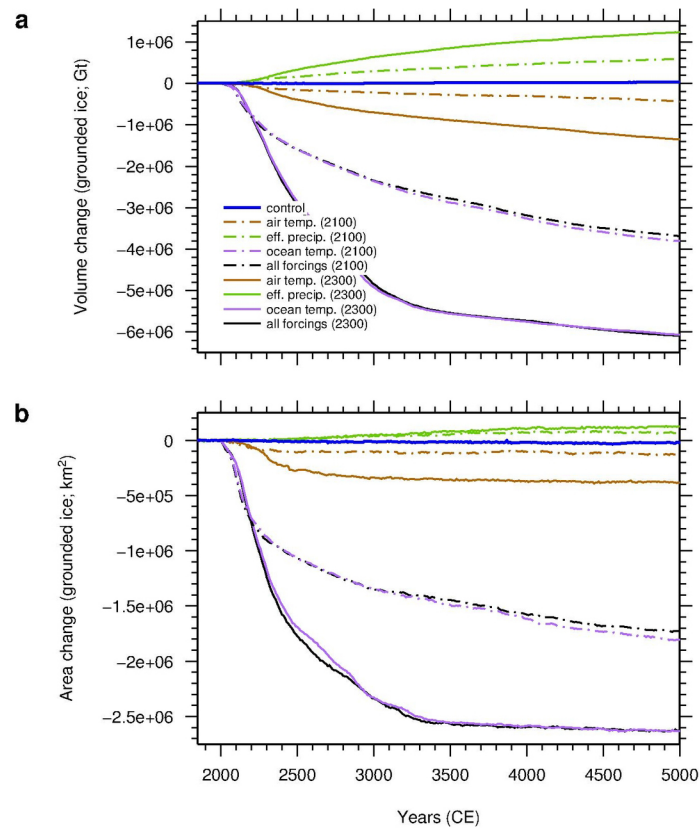
Antarctic-specific (60° – 90° S) projected temperature trends to 2300 CE based on CMIP5 values at 2100 CE and extended to 2300 CE following trajectories of global means from intermediate-complexity Earth system models^{3,47}. Precipitation and ocean temperature trends are calculated to follow those of

atmospheric temperatures, with magnitudes based on analysis of the CMIP5 data set indicating a 5.3% increase in precipitation per degree air temperature increase and a ratio of 0.25 for converting atmospheric to oceanic temperature changes. dT, change in air temperature.



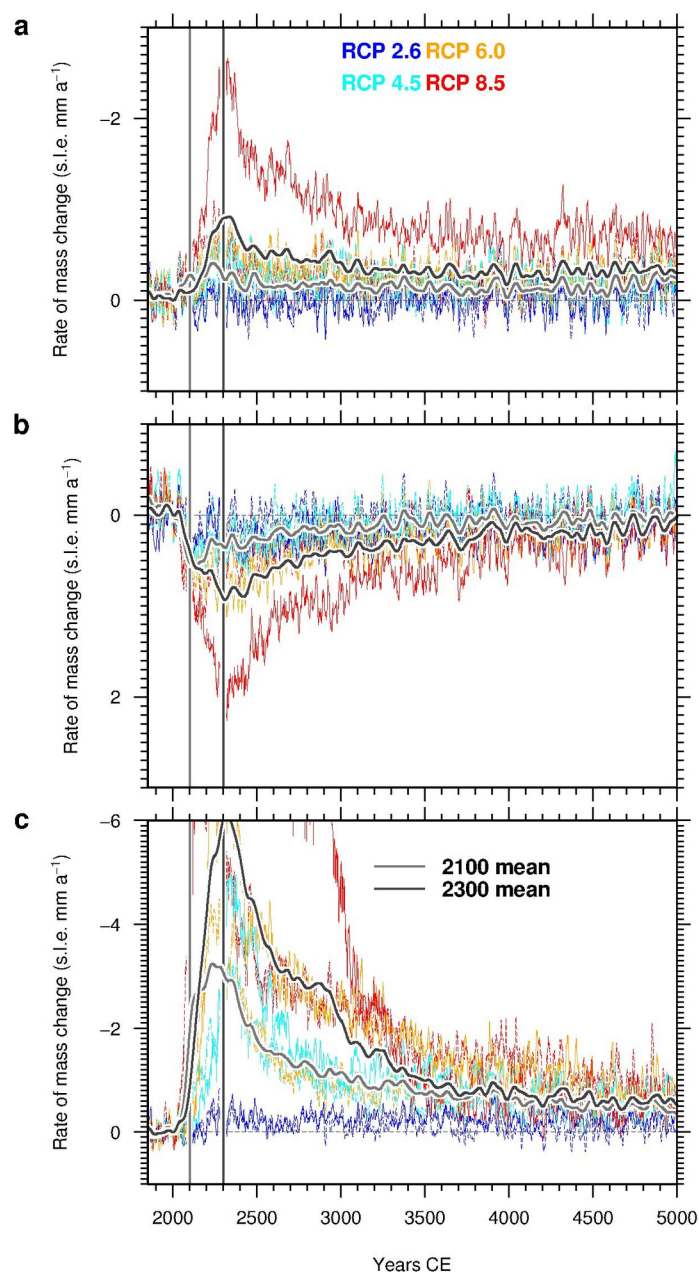
Extended Data Figure 3 | Model spin-up and fit to present-day. Ice-sheet geometry and surface velocities (ma^{-1} , metres per year) at the end of a 25,000-year evolutionary simulation. **a** and **b**, Observation-based⁵² (**a**) and modelled (**b**) ice-sheet extent and surface elevations. **c**, Comparison of ice

thicknesses shown in **a** and **b**. **d** and **e**, Measured⁵⁹ (**d**) and modelled (**e**) surface ice velocities. **f**, Comparison of the ice velocities shown in **d** and **e**. 'MeASURES' is the name of the published ice velocity dataset of ref. 59.



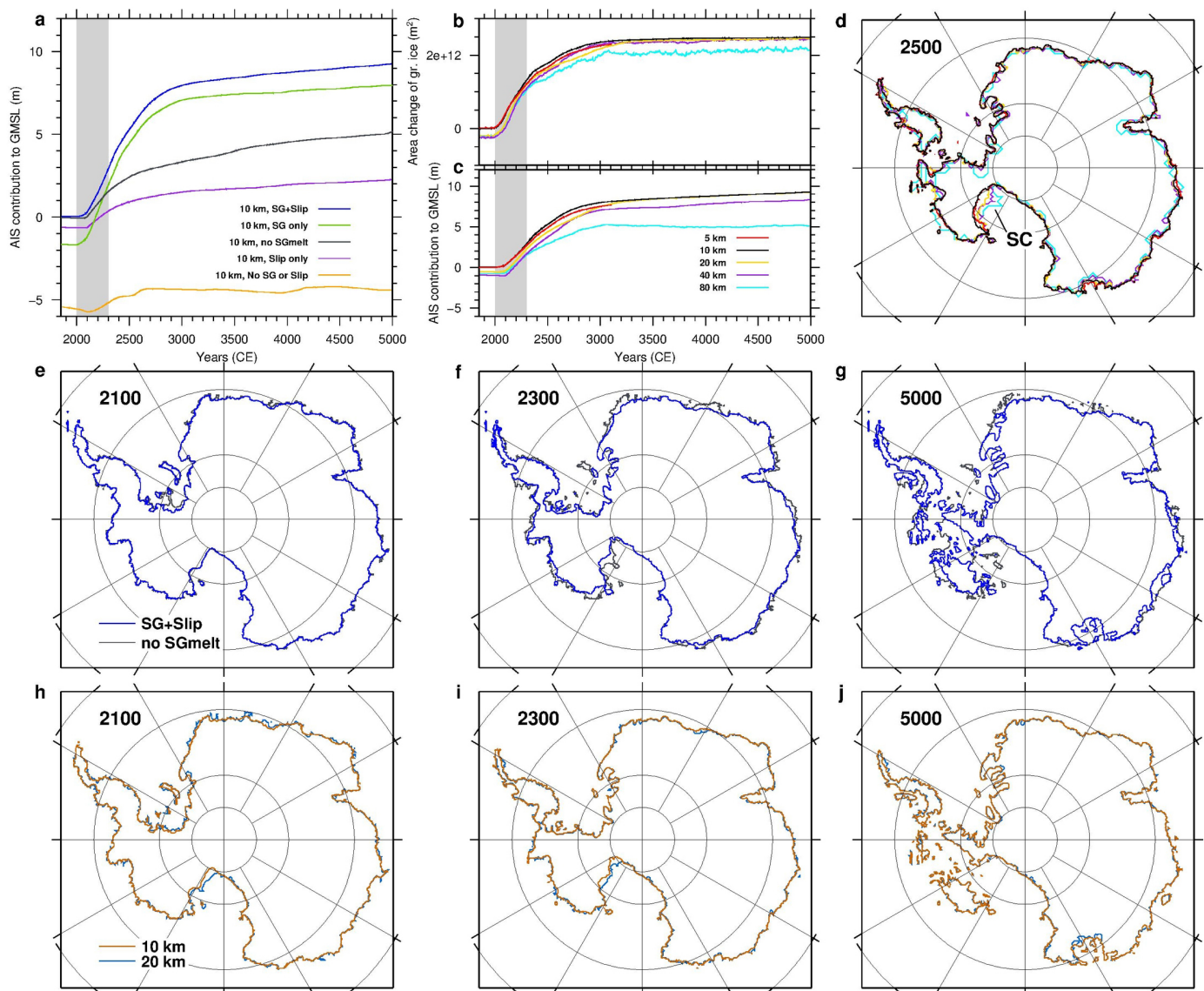
Extended Data Figure 4 | Multi-millennial changes in ice-sheet volume and area. Simulated changes in ice-sheet volume (**a**) and ice-sheet area (**b**) under single-parameter and combined forcings (ΔT_{air} , ΔP_{eff} and ΔSST), based on

simplified RCP scenarios for 100-year and 300-year forcing periods. Also shown in both panels is the control experiment (thick blue lines), illustrating little to no drift during the period of interest to 5000 CE.



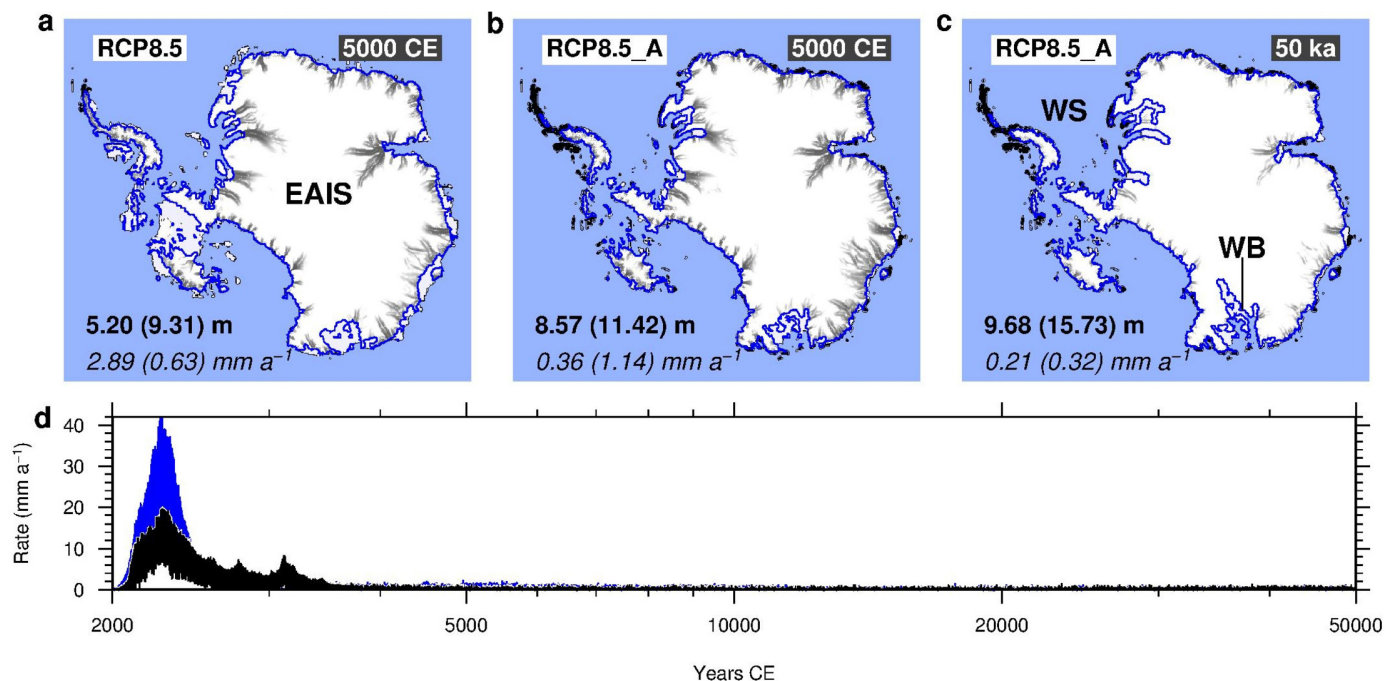
Extended Data Figure 5 | Multi-millennial changes in ice-sheet response to single-parameter environmental forcings. Bias-corrected rates of sea-level-equivalent ice-mass change (s.l.e. mm a⁻¹, millimetres of sea level equivalent per year) for each of the single-parameter simplified RCP forcing experiments. **a**, Rates of change under applied air temperature forcings peak at around 2240–2330 CE for both the 100-year and 300-year forcing experiments and decline thereafter, but both the 100-year and 300-year experiments exhibit rates that are still much larger than initial values by 5000 CE. **b**, Ice-mass rates of change as

in **a** but forced only with precipitation changes. Maxima for the 100-year and 300-year forcing experiments occur close to the end of the forcing period, reflecting little inherent lag. **c**, Rates of change in response to ocean forcing show much more elevated initial peaks compared to mass-loss rates in subsequent millennia. By the end of the run, rates of mass loss for both the 100-year and 300-year forcing experiments are still higher than at the beginning of the run. Data are shown relative to zero at 2000 CE.



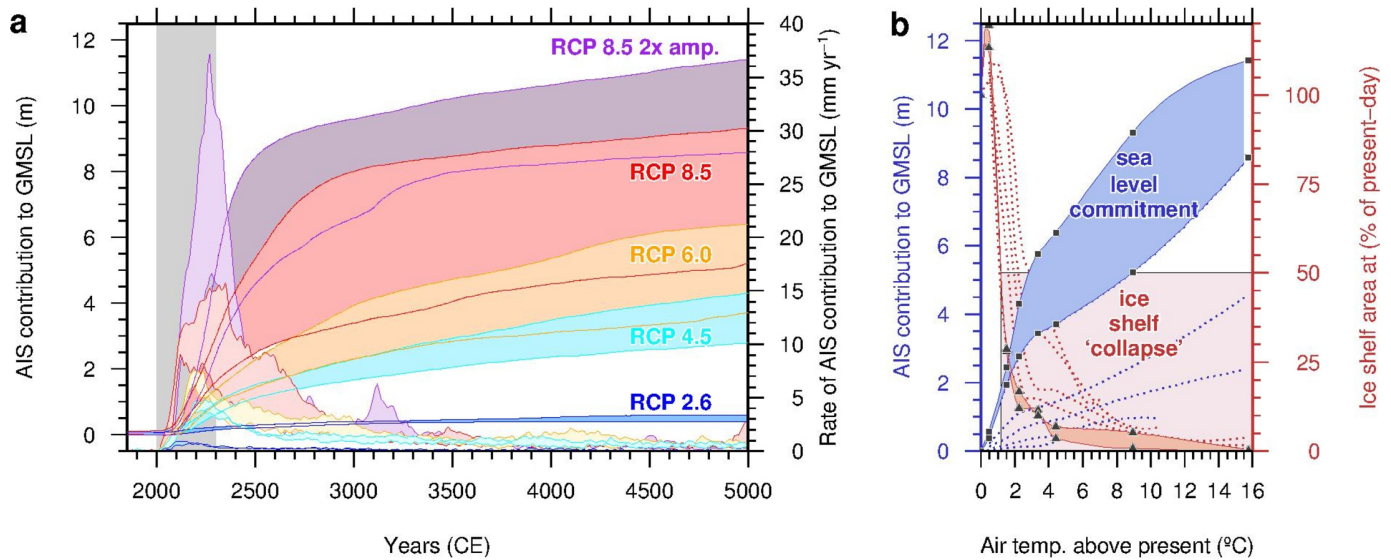
Extended Data Figure 6 | Model sensitivity and uncertainties. **a**, Modelled ice volume changes relative to the control run (in sea-level equivalent) for simulations in which only the grounding-line parameterization is altered. 'SG' and 'Slip' denote respectively the sub-grid and reduced traction grounding-line schemes employed in our simulations; 'no SGmelt' indicates an experiment in which only the sub-grid basal melt interpolation scheme is turned off (see Methods for details). Grey shading denotes period of applied forcing. **b** and **c**, Domain-integrated grounded (gr.) ice area (**b**) and sea-level-equivalent ice volume (**c**) trajectories under RCP 8.5 conditions for simulations in which the full sub-grid scheme is used and only the model resolution is changed. The 5-km simulation was not run beyond 3100 CE owing to the large computational

overhead. **d**, Grounding-line positions at 2500 CE for the experiments shown in **b** and **c**. The greatest differences occur in the Siple Coast area (SC). Note the very close agreement between the 10-km and 5-km simulations. **e–g**, Grounding-line locations under RCP 8.5 conditions at 2100 CE (**e**), 2300 CE (**f**) and 5000 CE (**g**) for experiments using the full grounding-line parameterization ('SG+Slip') compared to those in which the sub-grid basal melt interpolation is turned off ('no SGmelt'). **h–j**, Grounding-line locations under RCP 8.5 conditions at 2100 CE (**h**), 2300 CE (**i**) and 5000 CE (**j**) for 10-km and 20-km simulations that use the full grounding-line parameterization and resolution-specific stress balance tunings.



Extended Data Figure 7 | The effect of polar amplification. **a**, Geometry of the modelled Antarctic ice sheets under RCP 8.5 at 5000 CE, using both variants of the grounding-line scheme. Bold values and those in italics denote magnitudes and rates of sea-level contributions respectively. Leading values and those in parentheses relate to 'low' and 'high' scenarios respectively. Panels show ice extent for 'low' simulations; blue lines show grounding-line locations for 'high' simulations. Pale blue shading shows grounded ice lost in 'high' simulations but present in the 'low' scenario. **b**, Duplicate simulation to **a** but using a $2\times$ amplification of Antarctic temperatures beyond 2300 CE. Note the

greater sea-level contribution compared to **a**. **c**, The full equilibrium response of the polar amplification scenario shown in **b**. Note the greater loss of ice from the Wilkes Basin (WB) and eastern Weddell Sea (WS), resulting in a higher total sea-level contribution. Black areas denote ice-free land. **d**, Rate of ice loss for the $2\times$ amplification scenario for 'low' (black) and 'high' (blue) scenarios, illustrating that although the fastest contribution to sea level (2–4 m per century) occurs during the first millennium, slower mass loss continues for many millennia thereafter.



Extended Data Figure 8 | Antarctic contribution to GMSL. **a**, Predicted sea-level contribution from the AIS for 'high' and 'low' simulations (coloured lines) under each of the four RCP scenarios as well as one that includes $2\times$ amplification of Antarctic temperatures by 2300 CE (darker shading), based on coeval climatic and oceanic perturbations. The forced response (grey shading) represents 20% to 42% of the committed response by 5000 CE. Lighter shading between coloured lines shows rates of sea-level-equivalent ice loss for each scenario. **b**, Long-term sea-level commitment as a function of atmospheric warming (blue shading with squares). Intermediate response

curves for the 'low' simulations are shown in dotted lines. Red shading with triangles shows relationship between ice-shelf area and atmospheric warming for the near-equilibrium response and for intermediate stages (dotted lines). All curves in **b** are based on data from the four RCP scenario simulations, as well as one that includes $2\times$ amplification of Antarctic temperatures by 2300 CE, and two additional experiments whose maximum air temperature forcings are 1.5°C and 3.35°C . Pink shading defines the temperature range within which an ice-shelf extent less than 50% of present is simulated.

Extended Data Table 1 | CMIP5 multi-model ensemble mean environmental perturbations

		RCP 2.6	RCP 4.5	RCP 6.0	RCP 8.5
air temp. anomaly (°C)	Global	0.9	2.0	2.6	4.7
	Antarctic	0.8	1.6	2.3	4.2
precip. anomaly (%)	Global	4.2	10.9	13.2	22.6
	Antarctic	4.0	9.4	13.4	24.5
ocean temp. anomaly (°C)	Global	0.6	1.1	1.4	2.6
	Antarctic	0.2	0.4	0.5	1.0

Air temperature, precipitation and ocean temperature changes (ΔT_{air} , ΔP_{eff} and ΔSST) by 2100 for each of the RCP scenarios, expressed as perturbations from present.

Extended Data Table 2 | Parameter values found to best reproduce present-day Antarctic ice-sheet configuration

Parameter	Value	Units
Domain resolution (x, y)	20 <i>or</i> 10	km
Ice grid resolution (z)	0.024	km
Bedrock grid resolution (z)	0.1	km
Air temperature lapse rate	-0.008	°C m ⁻¹
SIA enhancement	1.2 <i>or</i> 1.47	-
SSA enhancement	0.5 <i>or</i> 0.57	-
Till porewater fraction	0.8	-
Eigen calving coefficient	5e+17	-
Thickness calving limit	220	m
Density of lithosphere	3300	kg m ⁻³
Flexural rigidity	5×10^{24}	Nm
Viscosity of mantle	1×10^{19}	Pa s

Note that different values of enhancement of shallow-ice and shallow-shelf approximations (SIA and SSA) are used for the 10-km- and 20-km-resolution simulations.

Inequality and visibility of wealth in experimental social networks

Akihiro Nishi^{1,2}, Hirokazu Shirado^{1,2}, David G. Rand^{1,3,4} & Nicholas A. Christakis^{1,2,5,6}

Humans prefer relatively equal distributions of resources^{1–5}, yet societies have varying degrees of economic inequality⁶. To investigate some of the possible determinants and consequences of inequality, here we perform experiments involving a networked public goods game in which subjects interact and gain or lose wealth. Subjects ($n = 1,462$) were randomly assigned to have higher or lower initial endowments, and were embedded within social networks with three levels of economic inequality (Gini coefficient = 0.0, 0.2, and 0.4). In addition, we manipulated the visibility of the wealth of network neighbours. We show that wealth visibility facilitates the downstream consequences of initial inequality—in initially more unequal situations, wealth visibility leads to greater inequality than when wealth is invisible. This result reflects a heterogeneous response to visibility in richer versus poorer subjects. We also find that making wealth visible has adverse welfare consequences, yielding lower levels of overall cooperation, inter-connectedness, and wealth. High initial levels of economic inequality alone, however, have relatively few deleterious welfare effects.

The unequal distribution of wealth in modern societies probably arose after we abandoned the relatively possession-free existence of hunter-gatherers^{7–9}, and it reflects several processes: individual variation in inborn traits (such as abilities, desires), differential access to environmental resources, and differential accumulation of wealth through transactions. Despite such inequality, humans have strong egalitarian preferences^{1–5}. What forces, then, lead to the emergence and maintenance of economic inequality? And what are the welfare implications of this inequality? We shed light on these questions using laboratory experiments that explore macro-level dynamics of economic inequality arising from micro-level cooperative interactions of individuals embedded within dynamic social networks^{10–12}. We focus on two dimensions: (1) initial conditions of wealth inequality (as a proxy for variation in initial endowments or private access to environmental resources), and (2) the local visibility of wealth.

We carried out a series of experiments with 1,462 subjects, divided among 80 sessions lasting an average of 30.0 minutes (s.d. = 7.13). Subjects were placed in groups with an average size of 17.21 (s.d. = 2.79) and arranged in a social network with an Erdős–Rényi random graph configuration in which 30% of ties were present (see Supplementary Information)^{10,11,13}; subjects were therefore initially connected to an average of 5.33 (s.d. = 0.98) neighbours. The subjects played a cooperation game lasting 10 rounds with their neighbours. In each round, all subjects chose whether to cooperate, by reducing their own wealth by 50 ‘units’ per neighbour in order to increase the wealth of all neighbours by 100 units each, or to defect by paying no cost and providing no benefits. Subjects made the same choice with all their neighbours. These interactions constituted the economic transactions, affecting each individual’s wealth and thus resulting in population-level changes in overall wealth and inequality. The arbitrary units were converted to real money at the end of the game (see Supplementary Information).

After making their cooperation choice, subjects were informed of the choices made by their neighbours. Then, subjects had the opportunity to change their neighbours by making or breaking ties. Specifically, 30% of all pairs of subjects were chosen at random in each round and given the opportunity to rewire their networks (this set-up was fixed across all manipulated conditions)^{10,11}. If a tie already existed between the two subjects, then one of the two was picked at random to be allowed to choose whether to voluntarily break the tie with the other; if a tie did not already exist between the two, both of them were given the option to form a tie and, if both approved, a new tie was formed. When making this decision, subjects were aware of whether the person to whom they might disconnect or connect had cooperated or defected in the past round. Thus, people could choose to alter a new subset of their social ties at each round; the network could be rewired; and subjects’ network degree (number of directly connected neighbours) and transitivity (the probability that any two of a focal subject’s neighbours are themselves connected) could change.

Within this basic setup, we then manipulated initial wealth inequality and wealth visibility (Extended Data Table 1 and Extended Data Figs 1 and 2). To manipulate initial wealth inequality, subjects were randomly assigned to one of three conditions. In the ‘no initial inequality’ condition, each subject started with the same initial endowment of 500 units. In the other two conditions, there was initial wealth inequality, such that ‘rich’ subjects received a larger initial endowment than ‘poor’ subjects. The endowments of the rich and poor were set to different levels of inequality such that the expected Gini coefficient (see Supplementary Information)¹⁴ at the beginning was either 0.0 (no initial inequality), 0.2 (low initial inequality), or 0.4 (high initial inequality). Importantly, the overall per capita initial wealth in all groups was equivalent (that is, 500 units); only the distribution of wealth varied. Subjects were randomly assigned to be rich or poor within the low and high initial inequality conditions, and they were randomly assigned to one of the nodes in the randomly generated network regardless of their endowment (see Fig. 1 for illustration, and also Supplementary Video 1).

Independent of baseline inequality, we also manipulated the visibility of local neighbours’ wealth. In the ‘invisible’ (private) condition, subjects only knew their own accumulated wealth. In the ‘visible’ condition, subjects could see their own accumulated wealth as well as the accumulated wealth of each of their directly connected neighbours. Subjects were informed whether each of their neighbours cooperated or not, regardless of the visibility condition of neighbours’ wealth. In both the visible and invisible set-ups, subjects had only local knowledge about their immediate neighbours and not global knowledge about the whole network.

Initial wealth inequality and visibility had joint and several effects on game dynamics. We begin by considering the persistence of wealth inequality (Fig. 2). Although the Gini coefficients in the invisible conditions converge at a low level (of roughly 0.16) by the later rounds, the Gini coefficients in the visible conditions vary persistently and depend on the initial level of inequality. We see a substantial interaction effect

¹Yale Institute for Network Science, Yale University, New Haven, Connecticut 06520, USA. ²Department of Sociology, Yale University, New Haven, Connecticut 06520, USA. ³Department of Psychology, Yale University, New Haven, Connecticut 06520, USA. ⁴Department of Economics, Yale University, New Haven, Connecticut 06520, USA. ⁵Department of Ecology and Evolutionary Biology, Yale University, New Haven, Connecticut 06520, USA. ⁶Department of Medicine, Yale University, New Haven, Connecticut 06520, USA.

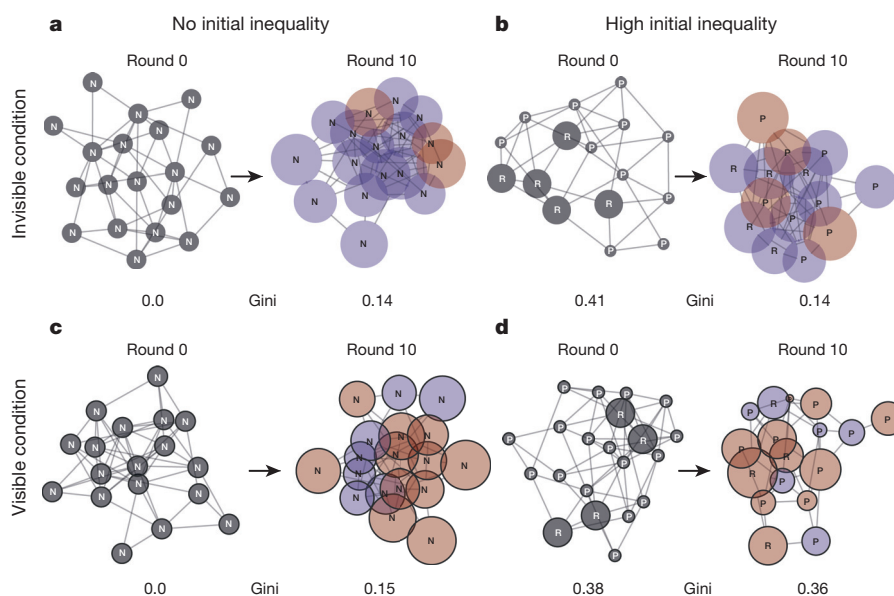


Figure 1 | The level of initial economic inequality and the visibility of connecting neighbours' wealth information are experimentally manipulated in dynamic human social networks. a–d, The states at round 0 (before interactions start) and at round 10 (after they end) in 4 out of 80 sessions ($n = 1,462$) are shown. The bold outline frame of a circular node indicates the 'visible' condition (wealth information is revealed to directly connected neighbours) and a non-bold outline frame indicates the 'invisible' condition (not revealed). Node size (area) indicates wealth (with bigger nodes

being richer). The letter in the node denotes the initial wealth to which subjects were randomly assigned: N is an initially non-poor/non-rich subject (in the no initial inequality condition), P is an initially poor subject, and R is an initially rich subject. Node colours represent the last move (purple, cooperate; red, defect; grey, no history). The Gini coefficient is also indicated (higher is more unequal). The Gini coefficient (for disposable income) is presently roughly 0.26 in Scandinavia and 0.39 in the United States. One of the three treatment conditions in our experiments (low initial inequality) is not shown.

between visibility and initial inequality on Gini over all rounds (two-way interaction $P = 0.043$; all P values determined using regression with standard errors clustered at the level of session and round; see Supplementary Table 2 and Supplementary Information for details). In the high initial inequality condition, making neighbours' wealth visible results in significantly higher levels of inequality compared to

when neighbours' wealth is invisible (difference in final round Gini = 0.104, $P = 0.004$) (Fig. 2, inset). In the low initial inequality condition, visibility again results in significantly higher inequality compared to the invisible condition, but to a lesser degree (difference in final round Gini = 0.0387, $P = 0.041$). Conversely, in the no initial inequality condition, visibility does not affect inequality (difference in final round Gini = 0.0185, $P = 0.450$). Thus, visibility serves to facilitate the persistence of whatever relative level of wealth inequality is initially present in the system, compared to what would have happened without visibility.

Examining groups of initially rich and poor subjects separately, we find that those individuals who are initially rich tend to be rich at the end, and, similarly, those who are initially poor tend to be poor at the end, regardless of whether the initial Gini coefficient is 0.2 or 0.4 (Extended Data Fig. 3). Although—in both the visible and invisible conditions—wealth distributions of initially rich and poor subjects gradually overlapped as the level of earned wealth increases in later rounds, few reversals of fortune occurred at the individual level (as also seen in labour markets¹⁵).

Turning to levels of average population wealth, we find that visibility has a substantial negative effect (Fig. 3a and Supplementary Table 4): despite the same payoffs and rules across conditions, overall wealth is significantly lower in the visible conditions compared to the invisible conditions (regression model coefficient = -489.6 , $P = 0.001$). The level of initial inequality is also negatively associated with overall wealth (coefficient = -669.6 , $P = 0.019$).

To further understand how visibility and inequality affect social welfare, we also examined cooperation and social tie formation. We find that the negative effect of visibility upon wealth accumulation is driven by a combination of two factors. First, cooperation rates are lower in the visible condition than the invisible condition (difference in cooperation frequency = -0.208 , $P < 0.001$; Fig. 3b and Supplementary Table 4), and do not differ based on the initial inequality—with a hypothetical change in the Gini coefficient from 0 to 1 being associated with a difference in cooperation frequency = -0.084 , $P = 0.445$. Second, there is lower social connectivity in the visible condition

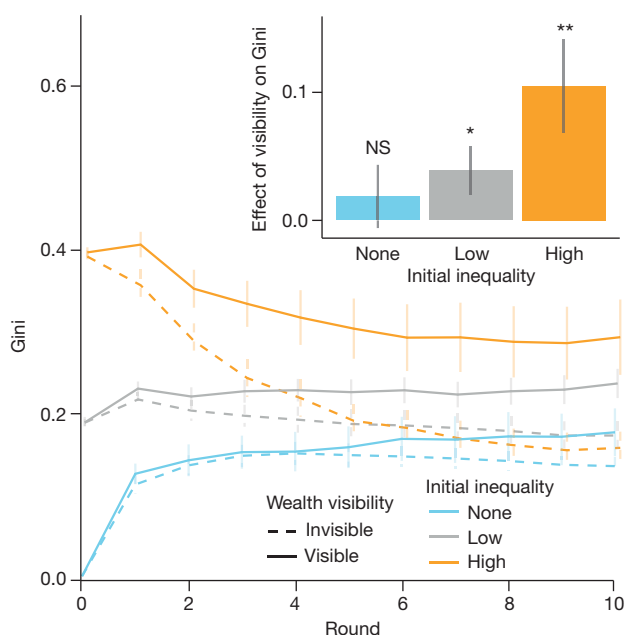


Figure 2 | Wealth visibility increases economic inequality (relative to invisibility) in the presence of initial inequality, but not in the presence of initial equality. The dynamics of the Gini coefficient in each of six settings (80 sessions total) is shown. Inset, the differences between the Gini coefficient over the ten rounds in the visible compared to the invisible condition (in the form of regression coefficients; see Supplementary Information). Error bars, mean \pm s.e.m. NS for $P \geq 0.05$, * $P < 0.05$, and ** $P < 0.01$.

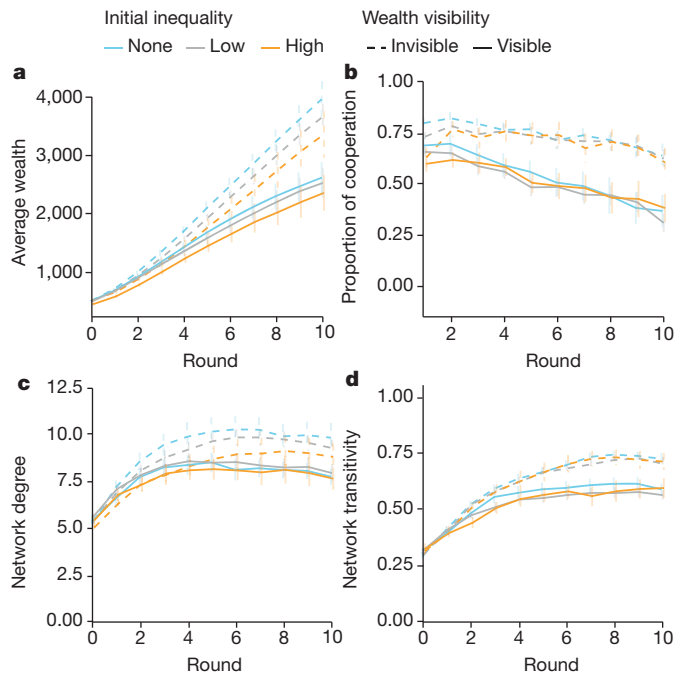


Figure 3 | Visibility of wealth undermines social welfare. **a–d**, Changes in average wealth (**a**), cooperation rate (**b**), network degree (number of connecting neighbours) (**c**), and network transitivity (probability of a focal subject's two neighbours being connected) (**d**), across rounds are shown (80 sessions total). Error bars, mean \pm s.e.m.

than in the invisible condition (difference in average degree = -0.991 , $P = 0.012$; Fig. 3c, Supplementary Table 4 and Extended Data Fig. 4), and there is no difference based on the initial inequality—with a hypothetical change in the Gini coefficient from 0 to 1 being associated with a difference in degree = -1.71 , $P = 0.167$. Visibility also seems to affect transitivity (difference in average transitivity = -0.096 , $P < 0.001$; Fig. 3d); however, after accounting for the rise in degree across rounds, and across treatments, neither visibility nor initial inequality affects transitivity (Extended Data Fig. 5b).

As average wealth (and overall group wealth) is roughly the multiplicative consequence of the cooperation rate and number of connecting neighbours, the dynamics of average wealth (Fig. 3a) can be explained by these dynamics of cooperation and degree. When we additionally explored these findings at the individual level, we found that subjects have a larger degree in the invisible condition as a consequence of there being a larger proportion of attractive neighbours (that is, cooperators at the last move) available in the social network (Supplementary Table 8). That is, visibility reduces cooperation which in turn reduces the appeal of social connections.

Although the dynamics of cooperation and degree can together explain the rise in average wealth, the macro-level dynamics of economic inequality that we observe in Fig. 2 require more micro-level analysis to fully explain. That is, the cooperation behaviour observed in Fig. 3b, when multiplied by the number of social connections shown in Fig. 3c, can explain the wealth shown in Fig. 3a, but it cannot explain the inequality shown in Fig. 2. Hence, to understand the dynamics of inequality, we examined how subjects exhibit different behavioural patterns of cooperation depending on their own wealth and on the average wealth of their neighbours, providing an individual-level understanding of the effect of visibility on inequality dynamics.

Figure 4 shows important heterogeneity in individual-level behaviours. When neighbours' wealth is visible, the level of initial inequality has a noticeable effect on how a subject's relative wealth affects the subject's cooperation (Fig. 4, right). In the high initial inequality condition, subjects who are locally and presently (that is, in the current round) richer than the average of their neighbours are less likely to

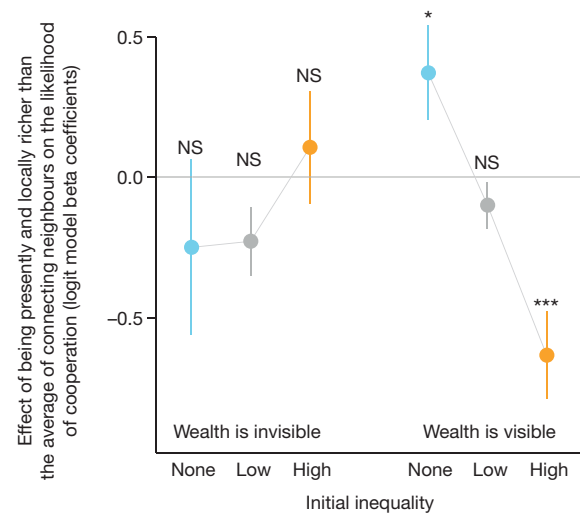


Figure 4 | Wealth visibility leads to exploitation under initial inequality, but fairness under initial equality. When wealth is visible (right), subjects richer than the average of their neighbours are more likely to cooperate in the 'no initial inequality' condition (blue, fairness scenario), but less likely to cooperate in the high initial inequality condition (orange, exploitation scenario). This behavioural pattern is not observed when connecting neighbours' wealth information is invisible (left). Shown are beta coefficients from logit models. P values indicate whether the coefficients are statistically different from 0.0. Error bars, point estimate \pm standard error. NS for $P \geq 0.05$, $*P < 0.05$, $***P < 0.001$. See Supplementary Information for details.

cooperate compared to those who are locally and presently poorer (regression model coefficient = -0.633 , $P < 0.001$). Moreover, we observe that this tendency is driven largely by richer-than-average subjects who defected in the prior round (coefficient = -0.997 , $P < 0.001$; Extended Data Fig. 6b); in an initially unequal world, we observe that defectors who are presently richer than connecting neighbours keep defecting and tend not to change their behaviour. This leads to an 'exploitation' scenario: poorer subjects are more likely to cooperate and invest in their local network, making them worse off relative to their neighbours and allowing the rich to get richer. As a result, richer subjects outperform poorer neighbours, leading to the increase in economic inequality (relative to the invisible condition) observed at the macro level.

Conversely, in the no initial inequality condition, we observe that subjects who are locally and presently richer than their neighbours are more likely to cooperate compared to poorer subjects (regression model coefficient = 0.370 , $P = 0.027$). Furthermore, this tendency is driven largely by richer subjects who cooperated in the prior round (coefficient = 0.805 , $P = 0.002$; Extended Data Fig. 6a); in an initially equal world, cooperators who are presently richer than connecting neighbours keep cooperating. This leads to a 'fairness' scenario in which the wealth of richer subjects is invested in their local network, allowing poorer neighbours to gain wealth. Thus, poorer subjects have the opportunity to catch up, and wealth visibility does not increase economic inequality relative to the invisible condition. Moreover, when considering mean difference in wealth, rather than Gini, this fairness behaviour leads to an actual reduction in inequality under visibility; see Extended Data Fig. 5a. As such effects are detected when connecting neighbours' wealth information is visible, but not when it is invisible (three-way interaction $P = 0.004$, Supplementary Table 7; $P > 0.05$ for all coefficients in the invisible condition in Fig. 4, left), these individual-level behaviours help to explain the macro-level results of our experiments. Moreover, agent-based simulations show that these patterns are sufficient to reproduce the observed inequality dynamics (see Supplementary Information and Extended Data Figs 7, 8 and 9).

In summary, we find that making wealth visible abets the persistence of experimentally induced inequality, compared to identical

circumstances where wealth is invisible. We also find that visibility has a corrosive overall effect on our laboratory 'societies', reducing overall cooperation, interconnectedness, and wealth. Thus, our experiments demonstrate that wealth visibility may be an important societal force, negatively affecting the dynamics of wealth and inequality, as well as social structure and cooperation. Surprisingly, our results are quite different with respect to the effect of initial wealth inequality. Rather than inequality being an 'enemy of cooperation', we find, in this setting, that inequality alone has relatively little effect on cooperation, interconnectedness or overall wealth accumulation. Thus, it is not inequality per se that is so problematic, but rather visibility that adversely affects cooperation here, regardless of what can be seen (that is, regardless of whether subjects are surrounded by an initially equal or unequal economic distribution).

Prior work regarding the role of inequality in contributions to public goods has reported mixed results^{16–20}, and the role of inequality in the evolution of cooperation has not been fully understood^{21–23}. Insofar as it is not inequality per se that affects cooperation in our experiments, but rather visibility, our results help shed light on these findings. Our results may also be relevant to norms in hunter-gatherer societies privileging less attachment to owned items⁸ and less ostentation²⁴. It is noteworthy that any (limited) wealth that is possessed in foraging societies is necessarily visible. Hence, it may not only be the surplus that arose with the agricultural revolution and fixed human settlements that contributed to inequality, but also the possibility of concealment that may be key. The mere ability to choose to conceal or display wealth might be relevant to how much inequality and cooperation arise in social groups.

Various psychological mechanisms may underlie the observed behaviours. For example, visibility may invoke neurological and psychological processes related to social comparison^{3,4,25,26}, and visibility may cause subjects to perceive the situation as a competition²⁷, to think that their wealth signals social position²⁸, or to fear being near last place²⁹, all of which might reduce cooperation. Our results are also consistent with findings regarding pay secrecy and worker productivity²⁶.

There are features potentially relevant to inequality that our experiments do not explore, for example: whether the initial resources are seen as earnings or windfalls; whether individuals producing public goods can earn more; whether the payoff structure, group size, network topology, or rewiring rate matter; or how peer sanctions or institutions (like taxation, courts or policing) affect the outcome. Another promising topic is the effect of allowing subjects to manipulate the visibility of wealth, in keeping with the theory of conspicuous consumption³⁰ and with notions of costly signalling. These are important directions for future work.

Although the results of laboratory experiments do not translate directly into the real world, the evidence presented here suggests that mechanisms that conceal personal wealth information might induce lower economic inequality, at least given an already high level of inequality. Given the widespread availability of wealth information as well as opportunities and desires to acquire and display wealth in contemporary societies, however, this would clearly not be easy to do. Conversely, when economic inequality is low, similarity could be more publicized, though this might sacrifice population-level economic growth.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 24 March; accepted 21 August 2015.

Published online 9 September 2015.

1. Dawes, C. T., Fowler, J. H., Johnson, T., McElreath, R. & Smirnov, O. Egalitarian motives in humans. *Nature* **446**, 794–796 (2007).
2. Fehr, E., Bernhard, H. & Rockenbach, B. Egalitarianism in young children. *Nature* **454**, 1079–1083 (2008).

3. Tricomi, E., Rangel, A., Camerer, C. F. & O'Doherty, J. P. Neural evidence for inequality-averse social preferences. *Nature* **463**, 1089–1109 (2010).
4. Dawes, C. T. *et al.* Neural basis of egalitarian behavior. *Proc. Natl Acad. Sci. USA* **109**, 6479–6483 (2012).
5. Engelmann, D. & Strobel, M. Inequality aversion, efficiency, and maximin preferences in simple distribution experiments. *Am. Econ. Rev.* **94**, 857–869 (2004).
6. Piketty, T. & Saez, E. Inequality in the long run. *Science* **344**, 838–843 (2014).
7. Apicella, C. L., Marlowe, F. W., Fowler, J. H. & Christakis, N. A. Social networks and cooperation in hunter-gatherers. *Nature* **481**, 497–501 (2012).
8. Apicella, C. L., Azevedo, E. M., Christakis, N. A. & Fowler, J. H. Evolutionary origins of the endowment effect: evidence from hunter-gatherers. *Am. Econ. Rev.* **104**, 1793–1805 (2014).
9. Smith, E. A. *et al.* Wealth transmission and inequality among hunter-gatherers. *Curr. Anthropol.* **51**, 19–34 (2010).
10. Rand, D. G., Arbesman, S. & Christakis, N. A. Dynamic social networks promote cooperation in experiments with humans. *Proc. Natl Acad. Sci. USA* **108**, 19193–19198 (2011).
11. Shirado, H., Fu, F., Fowler, J. H. & Christakis, N. A. Quality versus quantity of social ties in experimental cooperative networks. *Nature Commun.* **4**, 2814 (2013).
12. Rand, D. G., Nowak, M. A., Fowler, J. H. & Christakis, N. A. Static network structure can stabilize human cooperation. *Proc. Natl Acad. Sci. USA* **111**, 17093–17098 (2014).
13. Gracia-Lázaro, C. *et al.* Heterogeneous networks do not promote cooperation when humans play a Prisoner's Dilemma. *Proc. Natl Acad. Sci. USA* **109**, 12922–12926 (2012).
14. Allison, P. D. Measures of inequality. *Am. Sociol. Rev.* **43**, 865–880 (1978).
15. Calvó-Armengol, A. & Jackson, M. O. The effects of social networks on employment and inequality. *Am. Econ. Rev.* **94**, 426–454 (2004).
16. Cherry, T. L., Kroll, S. & Shogren, J. F. The impact of endowment heterogeneity and origin on public good contributions: Evidence from the lab. *J. Econ. Behav. Organ.* **57**, 357–365 (2005).
17. Chan, K. S., Mestelman, S., Moir, R. & Muller, R. A. Heterogeneity and the voluntary provision of public goods. *Exp. Econ.* **2**, 5–30 (1999).
18. Isaac, R. M. & Walker, J. M. Group-size effects in public-goods provision: the voluntary contributions mechanism. *Quart. J. Econ.* **103**, 179–199 (1988).
19. Sadrieh, A. & Verbon, H. A. A. Inequality, cooperation, and growth: an experimental study. *Eur. Econ. Rev.* **50**, 1197–1222 (2006).
20. Hackett, S., Schlager, E. & Walker, J. The role of communication in resolving common dilemmas: experimental evidence with heterogeneous appropriators. *J. Environ. Econ. Manage.* **27**, 99–126 (1994).
21. Abou Chakra, M. & Traulsen, A. Under high stakes and uncertainty the rich should lend the poor a helping hand. *J. Theor. Biol.* **341**, 123–130 (2014).
22. Wang, J., Fu, F. & Wang, L. Effects of heterogeneous wealth distribution on public cooperation with collective risk. *Phys. Rev. E* **82** (2010).
23. Kun, Á. & Diekmann, U. Resource heterogeneity can facilitate cooperation. *Nature Commun.* **4**, 2453 (2013).
24. Testart, A. The significance of food storage among hunter-gatherers: residence patterns, population-densities, and social inequalities. *Curr. Anthropol.* **23**, 523–537 (1982).
25. Gilbert, D. T., Giesler, R. B. & Morris, K. A. When comparisons arise. *J. Pers. Soc. Psychol.* **69**, 227–236 (1995).
26. Clark, A. E. & Oswald, A. J. Satisfaction and comparison income. *J. Public Econ.* **61**, 359–381 (1996).
27. Loughnan, S. *et al.* Economic inequality is linked to biased self-perception. *Psychol. Sci.* **22**, 1254–1258 (2011).
28. Moav, O. & Neeman, Z. Saving rates and poverty: the role of conspicuous consumption and human capital. *Econ. J.* **122**, 933–956 (2012).
29. Kuziemko, I., Buell, R. W., Reich, T. & Norton, M. I. "Last-place aversion": evidence and redistributive implications. *Quart. J. Econ.* **129**, 105–149 (2014).
30. Veblen, T. *The Theory of the Leisure Class: an Economic Study of Institutions* (Macmillan, 1899).

Supplementary Information is available in the online version of the paper.

Acknowledgements We thank C. Apicella, D. G. Alvarez, A. A. Arechar, D. Bergemann, G. Iosifidis, J. Jordan, J. H. Fowler, O. Moav, and A. Zaslavsky for comments. M. McKnight provided expert programming assistance and P. Treut provided technical support. The data reported in this paper are archived at Yale Institute for Network Science and are available upon request. A.N. was supported by the Japan Society for the Promotion of Science (JSPS). Support for this research was provided by a grant from the Robert Wood Johnson Foundation.

Author Contributions A.N., H.S., D.G.R., and N.A.C. designed the project. A.N. and H.S. conducted the experiments. A.N., H.S., and D.G.R. performed the statistical analyses. A.N., H.S., D.G.R. and N.A.C. analysed the findings. A.N., D.G.R., and N.A.C. wrote the manuscript.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to N.A.C. (nicholas.christakis@yale.edu).

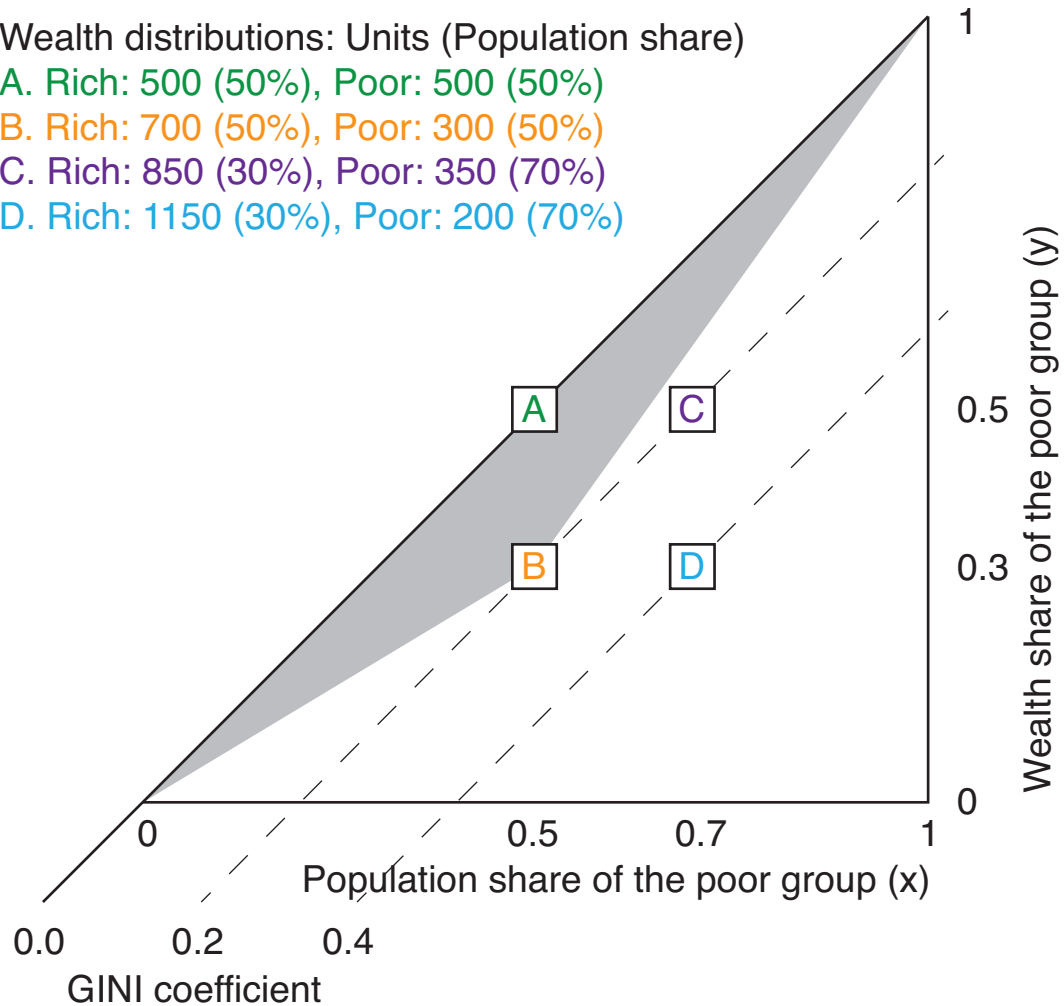
Wealth distributions: Units (Population share)

A. Rich: 500 (50%), Poor: 500 (50%)

B. Rich: 700 (50%), Poor: 300 (50%)

C. Rich: 850 (30%), Poor: 350 (70%)

D. Rich: 1150 (30%), Poor: 200 (70%)



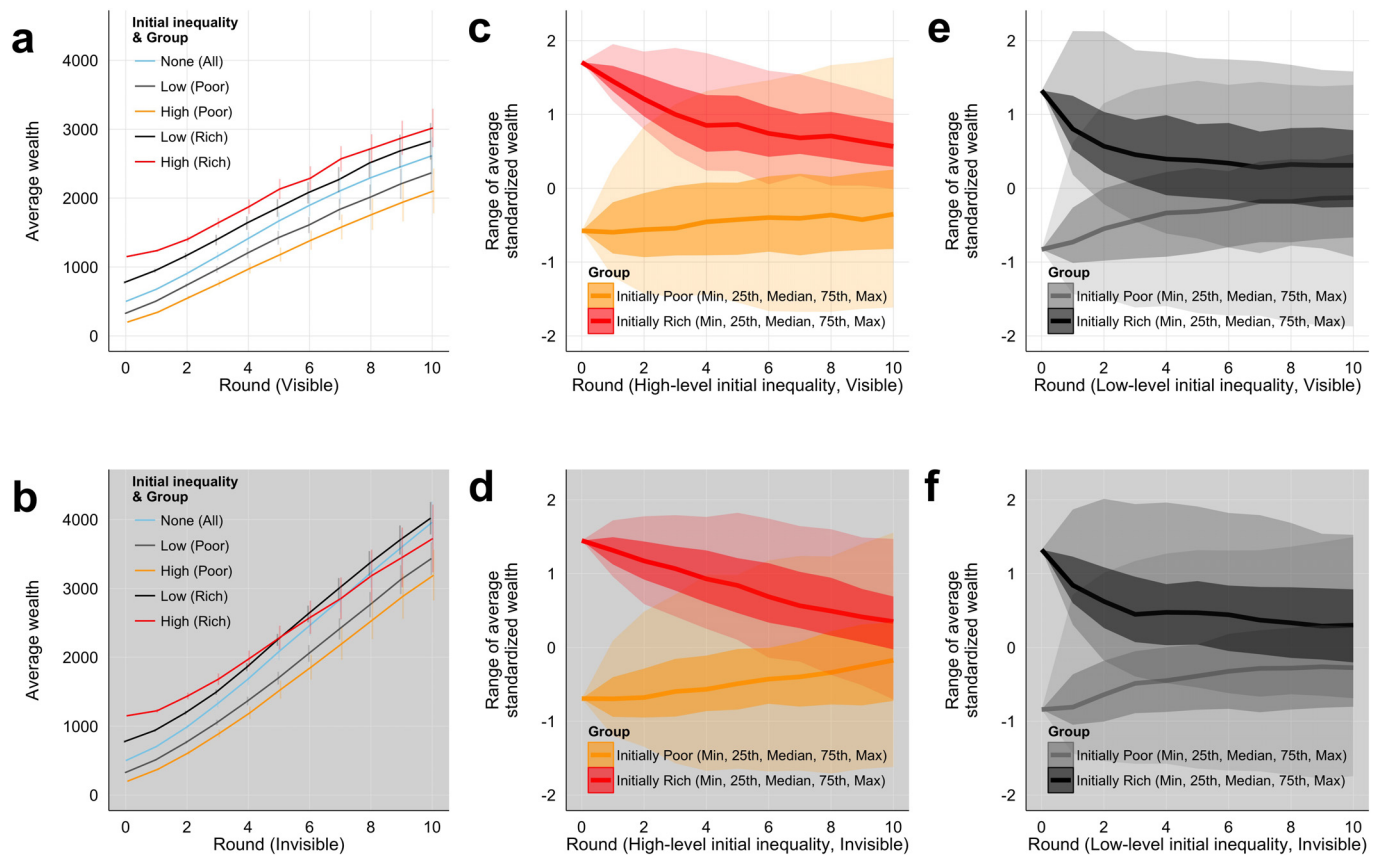
Extended Data Figure 1 | Lorenz curve description of the four wealth distributions. We prepared four different wealth conditions (A–D). For example, the shaded area for condition B divided by the area of the largest isosceles right triangle represents the Gini coefficient for condition B (that is, 0.2). Any points on the same dotted line achieve the same level of economic

inequality in terms of the Gini coefficient. Condition A is equivalent to any other condition on the line from (0,0) to (1,1). Conditions B and C are analysed together since they did not yield different analytical results (see Supplementary Information).

<div>Poorer Alter</div> <div>Richer Ego</div>	C	D
C	<div>+50</div> <div>GINI ↓</div> <div>Mean diff →</div> <div>+50</div>	<div>+100</div> <div>GINI ↓↓</div> <div>Mean diff ↓↓</div> <div>-50</div>
D	<div>-50</div> <div>GINI ↑↑</div> <div>Mean diff ↑↑</div> <div>+100</div>	<div>0</div> <div>GINI →</div> <div>Mean diff →</div> <div>0</div>

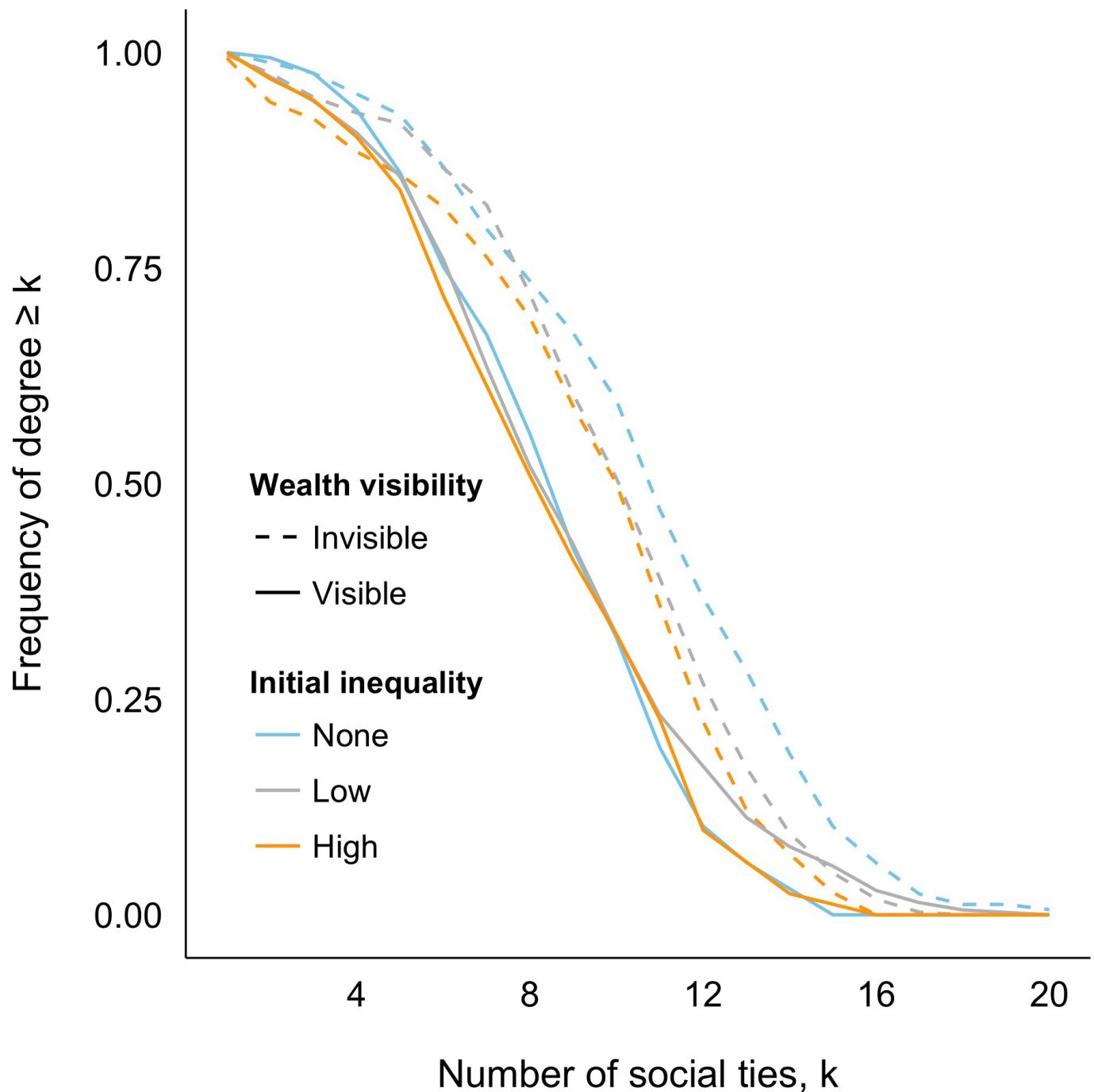
Extended Data Figure 2 | Rules in our experiments and the implied dynamics of the Gini coefficient and mean difference. The Gini coefficient is a relative measure of inequality, while the mean difference is an absolute measure of inequality. When we focus on a tie between two subjects (a richer ego and a poorer alter), there are four combinations in the choices of cooperation behaviours in a single round. For example, when a richer ego = C and a poorer alter = C (that is, when a richer ego cooperates and a poorer alter also cooperates), both of them pay 50 units, and obtain 100 units, in which case the payoff is +50 for both. As both of them get the same payoff, the mean difference between them does not increase or decrease (→). On the other hand,

the difference in wealth between them becomes less important in a relative manner, which leads to the reduction in the Gini coefficient between them (↓). The behaviours of the mean difference and Gini coefficient vary for the four combinations. C represents cooperation, and D represents defection. GINI represents local Gini coefficient of the focal two individuals, and ‘Mean diff’ represents the mean difference in wealth of the two subjects. GINI or mean difference can show the following outcomes: does not change (→), increases (↑), increases to a greater degree (↑↑), decreases (↓), or decreases to a greater degree (↓↓).



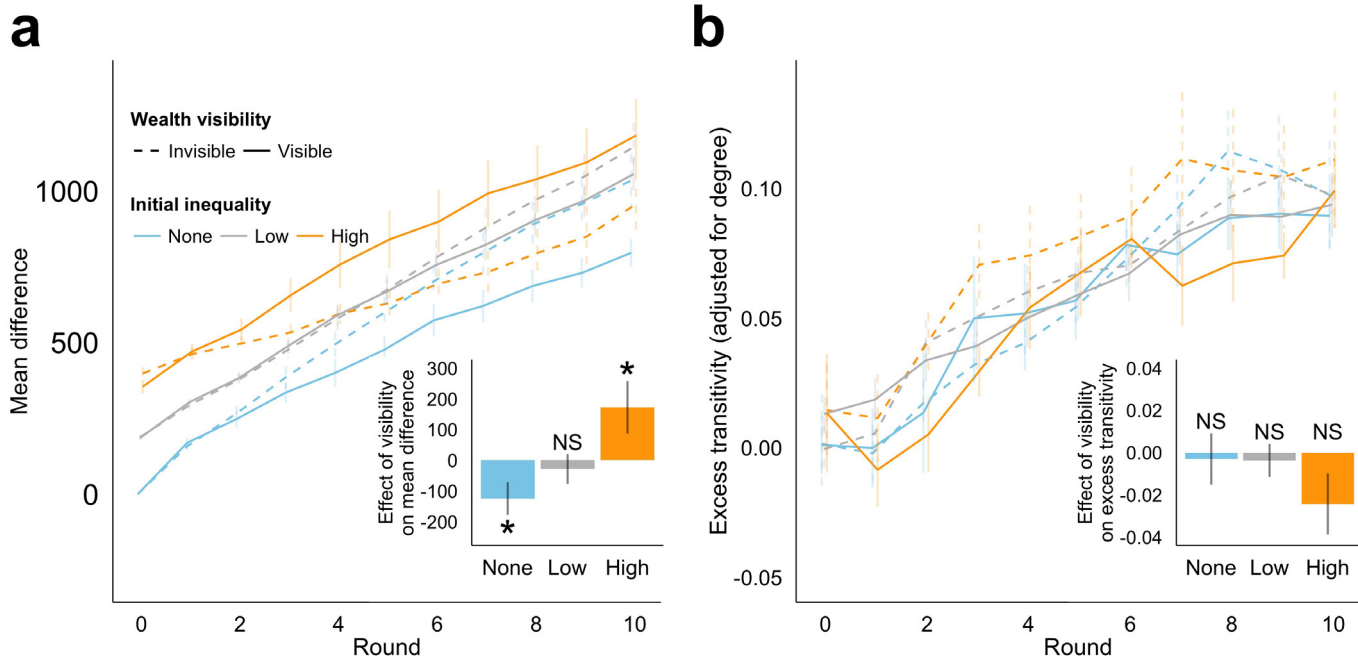
Extended Data Figure 3 | A majority of initially rich individuals stay wealthier than a majority of initially poor individuals over the ten rounds regardless of the initial conditions. **a, b**, The mean and standard error of mean (bar) of the average wealth of a group of initially rich individuals and of a group of initially poor individuals are calculated at each round for the visible condition (**a**) and for the invisible condition (**b**). Error bars, mean \pm s.e.m. **c–f**, For each round at each session, we standardized wealth of each individual (that is, (individual wealth–mean)/s.d.) and calculated the minimum (min), 25th percentile (25th), median, 75th percentile (75th), and maximum (max) of the standardized wealth of a group of initially rich individuals and a group of

initially poor individuals, separately. These figures show the trajectories of the mean of the five indicators (minimum, 25th, median, 75th, and maximum) among different sessions of the same initial condition (**c** for high-level initial inequality, visible; **d** for high-level initial inequality, invisible; **e** for low-level initial inequality, visible; and **f** for low-level initial inequality, invisible). Darker shades represent the mean of interquartile ranges (25th to 75th), lighter shades represent the mean of ranges (minimum to maximum), and the solid lines represent the mean of the median among the different sessions. Crossing of shades and medians between the two groups, if observed, implies the influence of the initial wealth difference on present wealth may be wiped out.



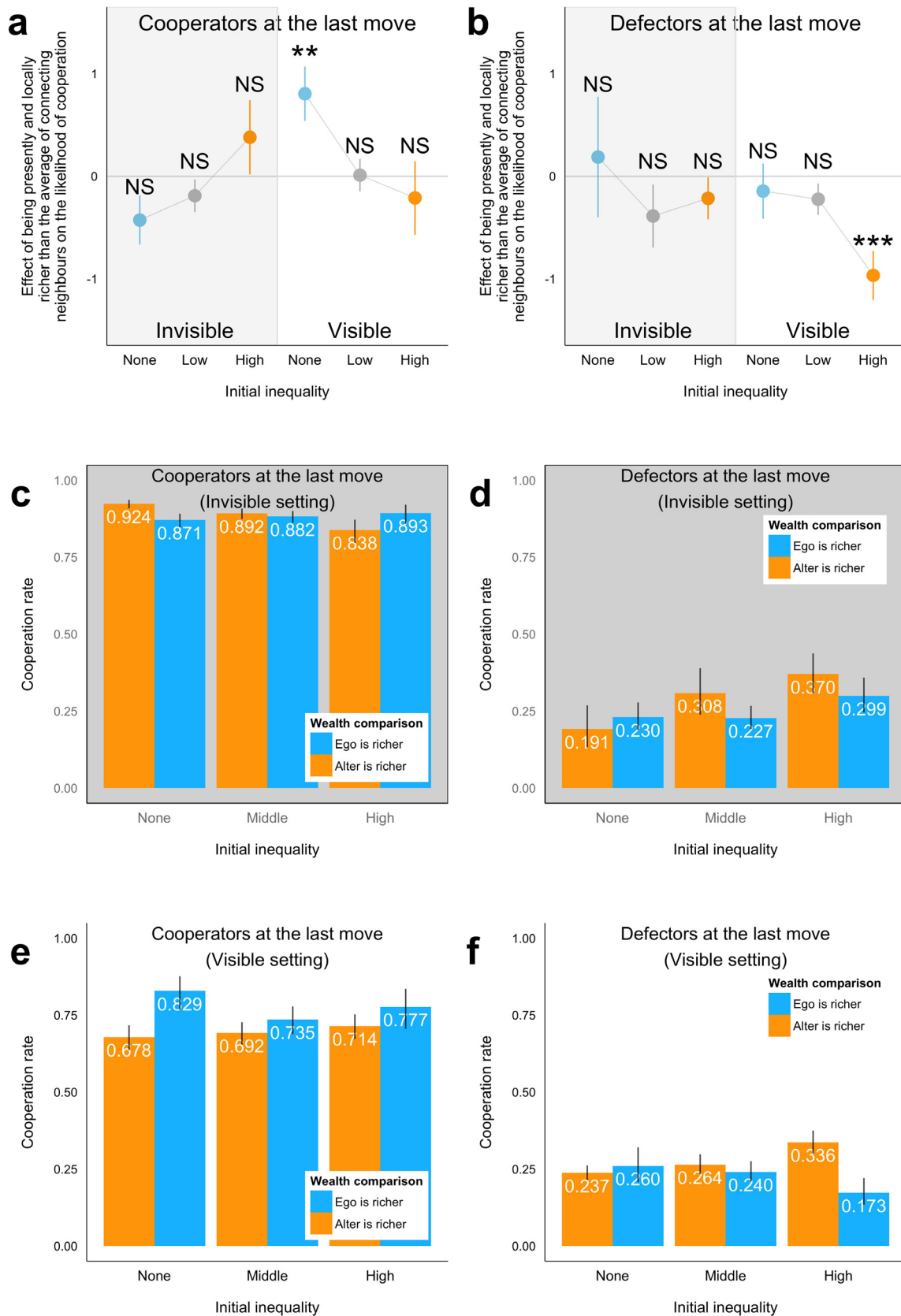
Extended Data Figure 4 | Cumulative degree distributions at the final (tenth) round. The proportion of subjects who have at least k social ties (degree) is calculated for each k (1 to 20) at each initial condition. Each distribution of the three initial inequality conditions in the invisible setting is significantly different from that in the visible setting (Kolmogorov–Smirnov test, $P < 0.01$), and has fatter tails. The pairwise comparison in different initial

inequality conditions at the same neighbours' wealth information setting (that is, none versus low, none versus high and low versus high) show those distributions are not significantly different (Kolmogorov–Smirnov test, $P > 0.12$) except none versus high in the invisible condition (Kolmogorov–Smirnov test, $P = 0.030$). The means of these distributions, by round, are shown in Fig. 3c.



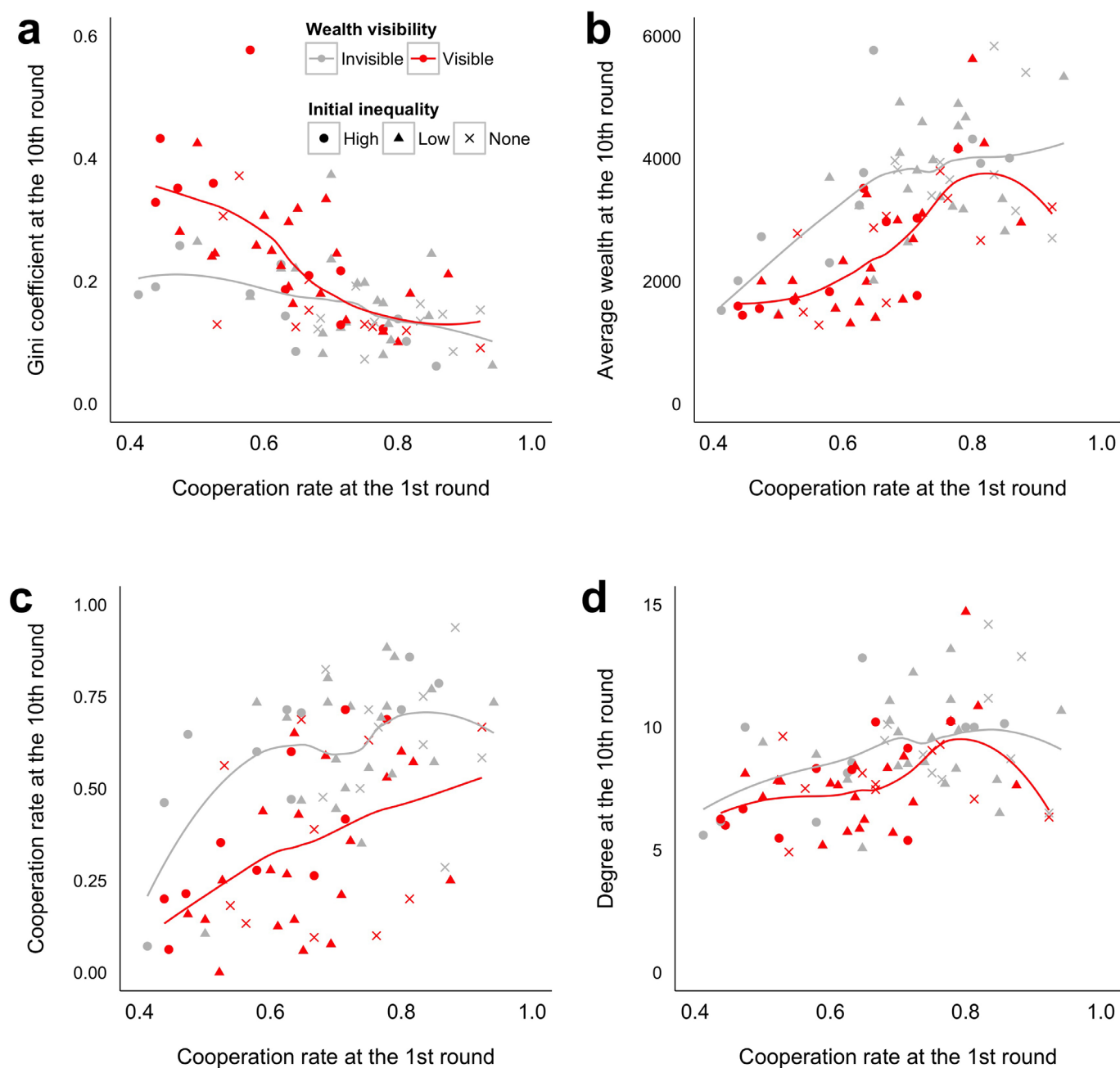
Extended Data Figure 5 | Changes in mean difference in wealth and in excess transitivity in the experimental conditions. **a**, The dynamics of mean difference in each of six settings is shown. Inset, the differences between mean difference at the first to tenth rounds in the visible compared to the invisible condition are shown separately for three different conditions of initial inequality (none, blue; low, grey; high, orange). Positive bars indicate that making neighbours' wealth visible increases mean difference in wealth. **b**, The dynamics of excess transitivity (transitivity adjusted for network degree at each session) in each of the six settings is shown. (See Fig. 3d for the dynamics of transitivity unadjusted for degree.) As a larger degree naturally results in a

larger transitivity, we calculate the expected value of transitivity given a certain network degree and a certain size in a random graph in simulations (10,000 iterations), and report the deviation of the observed transitivity from the expected transitivity (that is, observed transitivity minus expected transitivity). Inset, the differences between excess transitivity at the first to tenth rounds in the visible compared to the invisible condition are shown separately for three different conditions of initial inequality (none, blue; low, grey; high, orange). Negative bars indicate that making neighbours' wealth visible decreases excess transitivity. Error bars, mean \pm s.e.m. NS for $P \geq 0.05$, $*P < 0.05$.



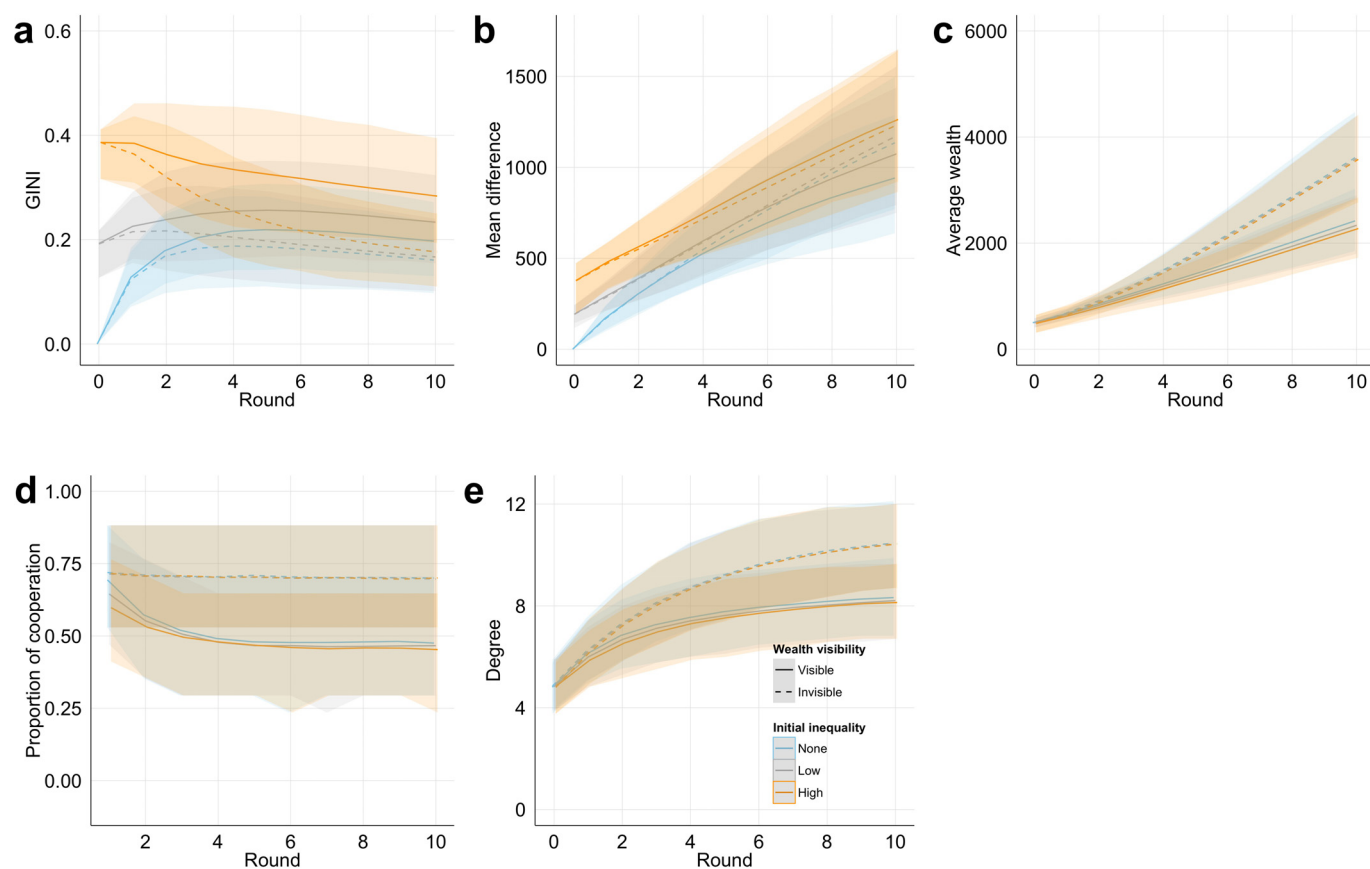
Extended Data Figure 6 | Additional results regarding behavioural mechanisms. **a, b**, Stratified results of Fig. 4 according to the prior move (cooperation or defection). Bars represent standard errors. **c–f**, Cooperation rate in the different conditions with respect to the variable of social comparison (ego (a focal individual) is richer or alters (given as the average of an ego's connecting neighbours) are richer) for each setting. For example, **e** shows that

richer subjects are more likely to cooperate (82.9%) when the initial economic inequality is set to none in the visible condition. Panel **f** shows that richer subjects are less likely to cooperate (17.3%) when the initial economic inequality is set to high in the visible condition. Error bars, mean \pm s.e.m. NS for $P \geq 0.05$, ** $P < 0.01$, *** $P < 0.001$.



Extended Data Figure 7 | Relationship between population-level cooperation rates at first round and outcomes at final rounds. a–d, Scatter plots of the first-round cooperation rate and the tenth-round Gini coefficient (a); average wealth (b); cooperation rate (c); and degree (d). Loess

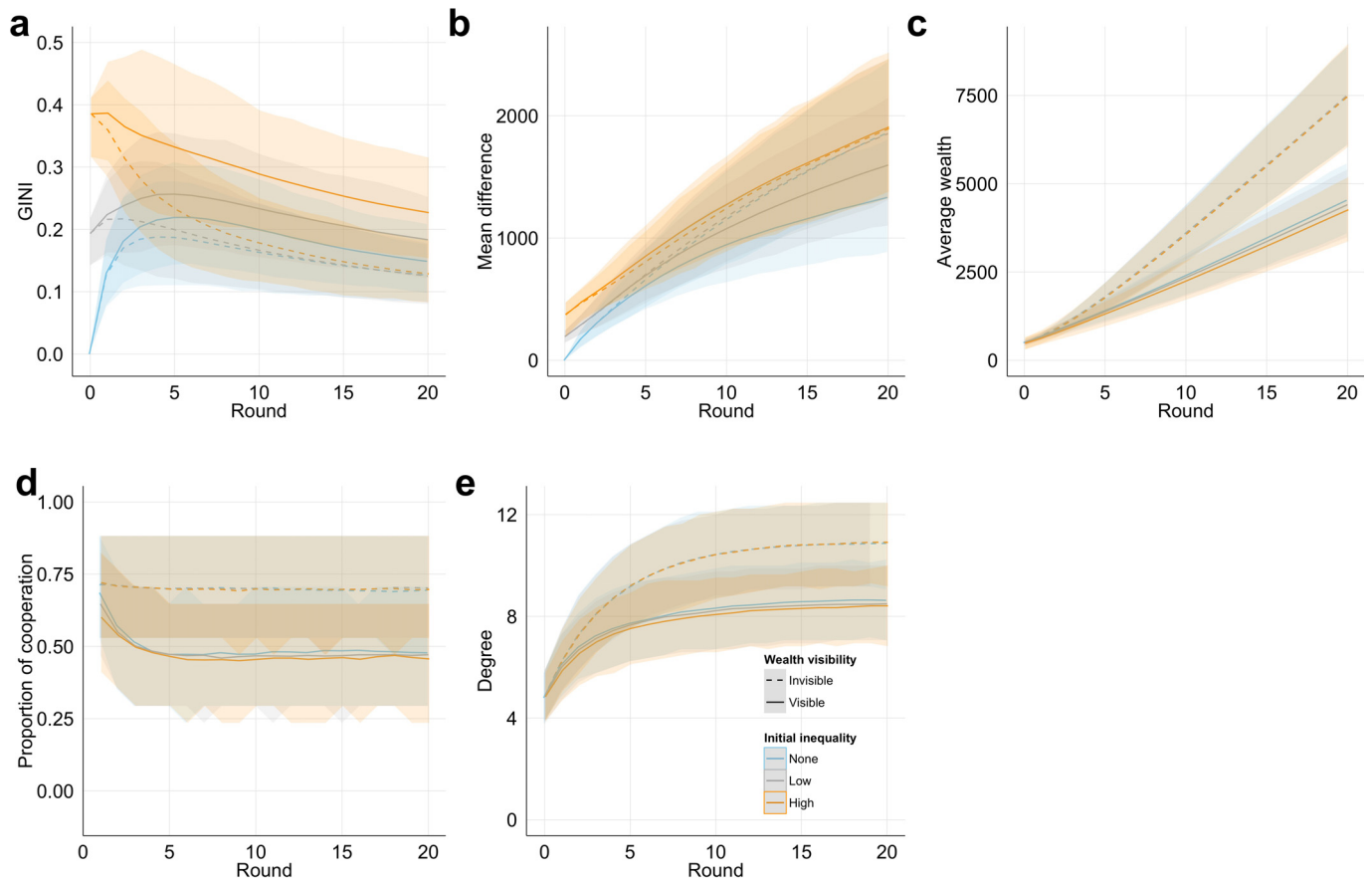
smoothed fitted curves are shown. The proportion of innately cooperative subjects in each session was not experimentally manipulated here (or in agent-based simulations).



Extended Data Figure 8 | Agent-based simulations reproduced the results that were observed in the online experiments with human subjects.

a–e, Results of agent-based simulations for Gini coefficient, mean difference,

average wealth, proportion of cooperation, and degree (interconnectedness) are shown, respectively. The medians (solid and dashed lines) and 90% confidence regions (shaded area, 5th percentile to 95th percentile) are presented.



Extended Data Figure 9 | Results of agent-based simulations with session up to 20 rounds show that the effect of visibility in Gini dynamics is robustly observed. **a–e**, The medians and 90% confidence regions (shaded area, 5th percentile to 95th percentile) are presented.

Extended Data Table 1 | Parameters for experimental settings

Visibility of connecting neighbours' wealth information	Wealth conditions	Initial wealth of rich individuals (%)	Initial wealth of poor individuals (%)	Level of initial economic inequality	Expected initial mean difference	Expected initial GINI	Number of games per session
Yes	A	500 (100%)		No	0	0.00	10
Yes	B	700 (50%)	300 (50%)	Low	200	0.20	10
Yes	C	850 (30%)	350 (70%)	Low	210	0.21	10
Yes	D	1150 (30%)	200 (70%)	High	399	0.41	10
No	A	500 (100%)		No	0	0.00	10
No	B	700 (50%)	300 (50%)	Low	200	0.20	10
No	C	850 (30%)	350 (70%)	Low	210	0.21	10
No	D	1150 (30%)	200 (70%)	High	399	0.41	10

Please refer to the Lorenz curve in Extended Data Fig. 1 for the wealth conditions A, B, C and D.

Forniceal deep brain stimulation rescues hippocampal memory in Rett syndrome mice

Shuang Hao^{1,2}, Bin Tang^{1,2}, Zhenyu Wu^{1,2}, Kerstin Ure^{1,3}, Yaling Sun^{1,3}, Huifang Tao^{1,3}, Yan Gao^{1,3}, Akash J. Patel^{1,4}, Daniel J. Curry⁴, Rodney C. Samaco^{1,3}, Huda Y. Zoghbi^{1,2,3,5,6,7} & Jianrong Tang^{1,2}

Deep brain stimulation (DBS) has improved the prospects for many individuals with diseases affecting motor control, and recently it has shown promise for improving cognitive function as well. Several studies in individuals with Alzheimer disease and in amnesic rats have demonstrated that DBS targeted to the fimbria–forix^{1–3}, the region that appears to regulate hippocampal activity, can mitigate defects in hippocampus-dependent memory^{3–5}. Despite these promising results, DBS has not been tested for its ability to improve cognition in any childhood intellectual disability disorder. Such disorders are a pressing concern: they affect as much as 3% of the population and involve hundreds of different genes. We proposed that stimulating the neural circuits that underlie learning and memory might provide a more promising route to treating these otherwise intractable disorders than seeking to adjust levels of one molecule at a time. We therefore studied the effects of forniceal DBS in a well-characterized mouse model of Rett syndrome (RTT), which is a leading cause of intellectual disability in females. Caused by mutations that impair the function of MeCP2 (ref. 6), RTT appears by the second year of life in humans, causing profound impairment in cognitive, motor and social skills, along with an array of neurological features⁷. RTT mice, which reproduce the broad phenotype of this disorder, also show clear deficits in hippocampus-dependent learning and memory and hippocampal synaptic plasticity^{8–11}. Here we show that forniceal DBS in RTT mice rescues contextual fear memory as well as spatial learning and memory. In parallel, forniceal DBS restores *in vivo* hippocampal long-term potentiation and hippocampal neurogenesis. These results indicate that forniceal DBS might mitigate cognitive dysfunction in RTT.

A deficit in contextual fear memory is one of the most reproducible and reliable outcome measures among RTT mouse models^{9–11}. Specifically, female *Mecp2*^{+/-} mice (hereafter referred to as RTT mice) have impaired contextual fear memory when tested 24 h after training¹¹. Because this deficit is readily quantifiable and accessible, we used fear memory as our first test of the effect of forniceal DBS in freely moving RTT mice. We implanted DBS electrodes in the fimbria–forix (FFx) of 6- to 8-week-old RTT mice and wild-type controls (Extended Data Fig. 1), guiding electrode placement with FFx-evoked potentials in the dentate gyrus (Fig. 1a–c). We divided mice into four groups after recovery: wild-type, sham; wild-type, DBS; RTT, sham; and RTT, DBS. Mice in both DBS groups received daily DBS treatment for 14 days while the two sham groups experienced the same procedures except for DBS. Based on widely used parameters for DBS in the clinic¹² and in rodents¹³, along with our own pilot testing, we set DBS at 130 Hz, 60 μ s pulse duration, and 1 h per day. Stimulus intensities were individually optimized to 80% of the threshold that elicits an instance of afterdischarge in the hippocampus^{3,14}. No seizures appeared under these DBS parameters. Three weeks after completing the two-week DBS protocol, we performed behavioural testing and subjected the

mice (now aged 14 weeks, Extended Data Fig. 1) to a fear conditioning paradigm to examine contextual fear memory and cued fear memory (see Methods).

Forniceal DBS significantly enhanced contextual fear memory in both wild-type (3 h, day 1 and day 3, $P < 0.05$) and RTT mice (3 h and day 1, $P < 0.05$; Fig. 1d). In fact, DBS restored contextual fear memory in RTT mice to wild-type levels: there was no difference between the DBS-treated RTT mice (3 h, $47.56 \pm 4.22\%$; day 1, $47.84 \pm 4.16\%$) and sham-treated wild-type mice (3 h, $44.87 \pm 3.60\%$; day 1, $45.97 \pm 3.69\%$). Notably, forniceal DBS did not alter cued fear memory (Fig. 1e), even though the FFx also projects to the amygdala¹⁵. All the mice that received DBS/sham treatment responded to tone presentation (Extended Data Fig. 2e–h), but less than the animals that were implanted and did not experience the 2-week DBS/sham procedures (Extended Data Fig. 2b–d). Further analysis indicated that the longer period of handling and exposure (for example, daily transportation, connection/disconnection of the wires, and staying in the DBS/sham chamber for 1 h per day) increased the motor activity and decreased the anxiety levels in DBS/sham-treated mice (Extended Data Fig. 3). These changes likely reduced the fear responses to the tone, and the conditioning context, in general (Fig. 1d, e and Extended Data Fig. 2a, b).

Forniceal DBS did not improve levels of locomotion, anxiety, pain threshold or motor learning (Extended Data Figs 3 and 4a, b) as well as motor coordination, social behaviour and body weight in RTT mice, although there were differences between RTT mice and wild-type controls in these features (Extended Data Figs 4c–f and 5a, b). Forniceal DBS thus specifically rescued contextual memory impairment in RTT mice without evident off-target effects.

To determine whether forniceal DBS would improve spatial cognition, which is also hippocampus-dependent, we trained new cohorts of mice, who received the same DBS/sham procedures, in a hidden platform version of the water maze task¹⁶ (Extended Data Fig. 1). Sham-treated RTT mice needed more time than sham-treated wild-type mice to locate the hidden platform across the training trials, spent less time in the target quadrant, and had fewer platform area crossings in the probe test (Fig. 2a). In wild-type mice, DBS significantly enhanced spatial learning compared to the sham group (Fig. 2b). Treatment made no difference during the probe test, probably because of a ceiling effect in sham-treated wild-type animals. We observed an even stronger effect of DBS in RTT mice: forniceal DBS enhanced not only spatial learning but also spatial memory retrieval (Fig. 2c). Again, the rescue was so strong that there was no difference between DBS-treated RTT and sham-treated wild-type groups in latencies to the hidden platform, time in target quadrant, or platform area crossings (Fig. 2d). Visible platform training confirmed that neither MeCP2 level nor forniceal DBS altered visual or sensorimotor skills (Extended Data Fig. 5c–e).

¹Jan and Dan Duncan Neurological Research Institute, Texas Children's Hospital, Houston, Texas 77030, USA. ²Department of Pediatrics, Baylor College of Medicine, Houston, Texas 77030, USA.

³Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas 77030, USA. ⁴Department of Neurosurgery, Baylor College of Medicine, Houston, Texas 77030, USA. ⁵Program in Developmental Biology, Baylor College of Medicine, Houston, Texas 77030, USA. ⁶Department of Neuroscience, Baylor College of Medicine, Houston, Texas 77030, USA. ⁷Howard Hughes Medical Institute, Baylor College of Medicine, Houston, Texas 77030, USA.

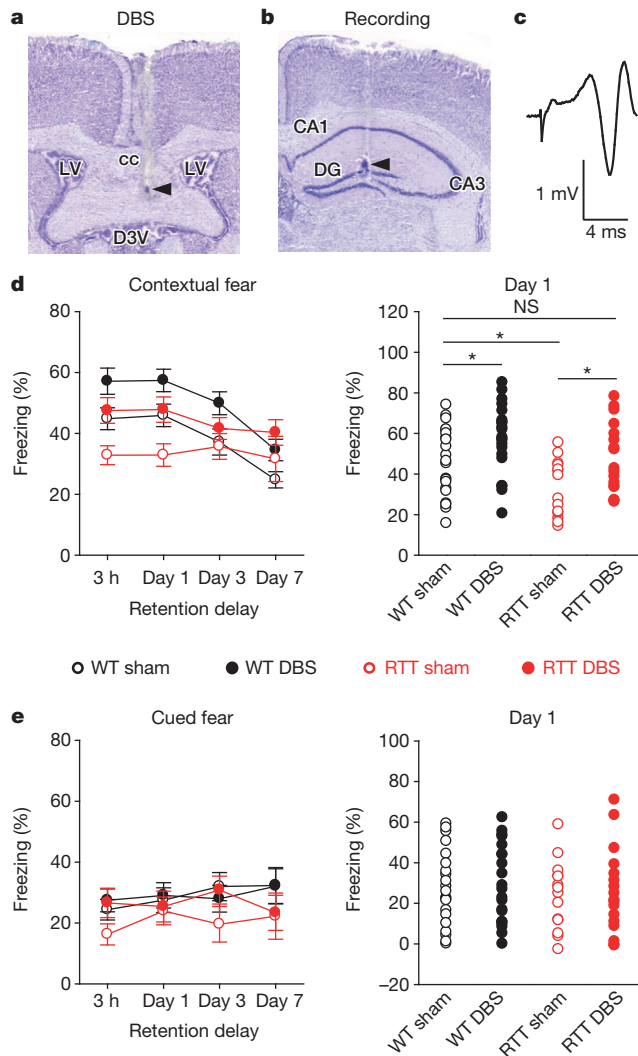


Figure 1 | Forniceal DBS restores contextual fear memory in RTT mice.

a, b, Photomicrographs illustrating DBS electrode placement (arrowheads) in the FFx (**a**) and the recording electrode in the dentate gyrus (**b**). cc, corpus callosum; LV, lateral ventricle; D3V, dorsal third ventricle; DG, dentate gyrus. **c**, Representative evoked potential trace of the FFx pathway recorded in the dentate. **d**, Forniceal DBS enhanced contextual fear memory in both wild-type (WT) and RTT mice (DBS-treated wild-type mice, $n = 21$; sham-treated wild-type mice, $n = 21$; DBS-treated RTT mice, $n = 17$; sham-treated RTT mice, $n = 14$). There were significant main effects on freezing time among the four groups (two-way repeated-measures ANOVA: group, $F_{3,69} = 5.67$, $P = 0.002$; day, $F_{3,180} = 6.44$, $P < 0.001$; group \times day interaction, $F_{9,180} = 2.15$, $P = 0.027$). Within-genotype analysis revealed a significant DBS effect in both wild-type ($F_{1,40} = 8.50$, $P = 0.006$) and RTT mice ($F_{1,29} = 6.44$, $P = 0.016$). DBS in RTT mice restores contextual fear memory to wild-type levels (DBS-treated RTT mice vs. sham-treated wild-type mice: group, $F_{1,36} = 2.76$, $P = 0.105$). Comparison of contextual fear memory on day 1 among the four groups revealed a significant main effect (two-way ANOVA followed by Tukey's post hoc test: genotype, $F_{1,69} = 8.39$, $P = 0.005$; treatment, $F_{1,69} = 11.41$, $P = 0.001$). **e**, Cued fear memory of mice tested in **d**. There was no difference in cued fear memory between groups over any time point (main effect: group, $F_{3,69} = 0.88$, $P = 0.456$; day, $F_{3,180} = 1.65$, $P = 0.179$; group \times day interaction, $F_{9,180} = 0.89$, $P = 0.538$) or on day 1 (genotype, $F_{1,69} = 0.64$, $P = 0.428$; treatment, $F_{1,69} = 0.11$, $P = 0.741$). * $P < 0.05$; NS, not significant. Data presented as mean \pm s.e.m. Scattergrams show individual values.

Since the RTT mice used in this study are impaired in both hippocampus-dependent memory¹¹ and *in vitro* hippocampal long-term potentiation (LTP)⁸, they provide an ideal setting in which to examine whether DBS alters synaptic plasticity. We implanted RTT and wild-type mice with DBS electrodes in the FFx, stimulation electrodes in

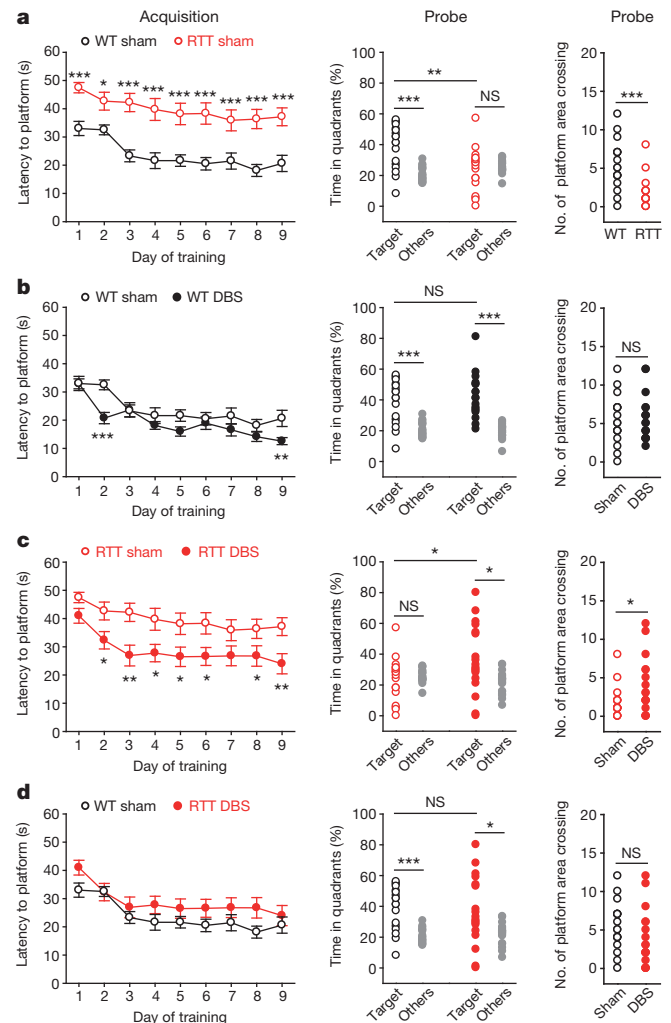


Figure 2 | Forniceal DBS rescues spatial learning and memory in RTT mice.

In the water maze task, all mice were trained with a hidden platform for 9 days followed by a probe test without the platform 24 h after the last training. There were significant main effects of escape latencies among the four groups ($n = 18$ mice per group) during acquisition training (two-way repeated-measures ANOVA: group, $F_{3,68} = 20.74$, $P < 0.001$; day, $F_{8,544} = 19.72$, $P < 0.001$). **a**, Sham-treated RTT mice showed increased escape latencies during training (genotype, $F_{1,34} = 35.30$, $P < 0.001$; day, $F_{8,272} = 7.06$, $P < 0.001$), but decreased time in target quadrant ($P < 0.01$) and fewer platform area crossings during the probe test ($P < 0.001$) than sham-treated wild-type controls. **b**, Forniceal DBS decreased escape latencies during training in DBS-treated wild-type mice compared to sham-treated wild-type controls (treatment, $F_{1,34} = 5.94$, $P = 0.020$; day, $F_{8,272} = 17.10$, $P < 0.001$; treatment \times day, $F_{8,272} = 2.19$, $P = 0.028$). There was no difference in time spent in the target quadrant ($P > 0.05$) or in the number of platform area crossings ($P > 0.05$) between DBS and sham groups during the probe test. **c**, The DBS-treated RTT mice showed shorter escape latencies during training (treatment, $F_{1,34} = 10.31$, $P = 0.003$; day, $F_{8,272} = 6.13$, $P < 0.001$) but more time in the target quadrant ($P < 0.05$) and platform area crossings ($P < 0.05$) during the probe test than sham-treated RTT controls. **d**, There was no difference between DBS-treated RTT mice and sham-treated wild-type controls in escape latencies during training (group, $F_{1,34} = 2.91$, $P = 0.097$; group \times day interaction, $F_{8,272} = 0.80$, $P = 0.606$), time in the target quadrant ($P > 0.05$), or number of crossings of the platform area ($P > 0.05$) during the probe test. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$; NS, not significant (Tukey's post hoc test during acquisition; two-tailed unpaired *t*-test between groups and two-tailed paired *t*-test within group during probe). Data presented as mean \pm s.e.m. Scattergrams show individual values.

the perforant path for the LTP test, and a recording electrode in the dentate gyrus (using evoked potentials as a guide). As with the behavioural studies described above, the mice underwent 2 weeks of DBS

followed by a 3-week interval before LTP testing (Extended Data Fig. 1). These DBS/sham procedures did not alter the hippocampal neural excitability (Extended Data Fig. 6a). LTP was induced on day 0 and monitored across 5 days after induction in awake, freely moving mice. We quantified the population spike amplitude in each of the four groups of animals¹⁷ (Fig. 3a). Neither DBS nor sham treatment altered baseline synaptic transmission in sham-treated wild-type, sham-treated RTT or DBS-treated wild-type mice, but DBS slightly reduced the magnitude of the evoked responses in DBS-treated RTT animals (Extended Data Fig. 6b). There was no difference in the stimulus intensities used for LTP induction among the four groups (wild-type, sham, $83.91 \pm 13.29 \mu\text{A}$; wild type, DBS, $69.50 \pm 7.58 \mu\text{A}$; RTT, sham, $73.50 \pm 14.58 \mu\text{A}$; RTT, DBS, $68.62 \pm 13.87 \mu\text{A}$; $P > 0.05$). As expected, sham-treated RTT mice showed impaired *in vivo* LTP (1 h after induction, $160.37 \pm 14.08\%$) compared to sham-treated wild-type controls ($256.06 \pm 27.50\%$) (Fig. 3b). Forniceal DBS, however, enhanced LTP in both wild-type mice (Fig. 3c; DBS, $381.75 \pm 26.13\%$; sham, $256.06 \pm 27.50\%$) and RTT animals (Fig. 3d; DBS, $254.93 \pm 22.96\%$; sham, $160.37 \pm 14.08\%$) to the degree that there was no difference between DBS-treated RTT and sham-treated wild-type groups (Fig. 3e). Forniceal DBS thus restored hippocampal LTP in the perforant path/dentate pathway in RTT mice.

Because hippocampal neurogenesis contributes to hippocampal LTP and hippocampus-dependent memory^{18–20}, and because DBS in other afferent pathways of the hippocampus increases dentate neurogenesis^{21–23}, we explored whether forniceal DBS might exert its effects

by stimulating hippocampal neurogenesis. We first observed that two-hour unilateral forniceal DBS stimulated the activity of dentate neurons, as indicated by increased expression of the immediate early gene *Fos* (Extended Data Fig. 7a). Each day after DBS over the 2 weeks of DBS/sham treatment, we injected RTT and wild-type mice with 5-bromo-2'-deoxyuridine (BrdU) to mark newborn cells. We quantified dentate neurogenesis by the numbers of cells positive for BrdU, DCX (doublecortin, to label the immature neurons), or double labelled in each of the four groups (Fig. 4a–d). We found that baseline levels of dentate neurogenesis in sham-treated RTT mice (BrdU, 571 ± 99.84 ; DCX, 914 ± 242.76) were significantly lower than in sham-treated wild-type controls (BrdU, $1,218 \pm 175.57$; DCX, $2,059 \pm 381.49$) (Fig. 4e, f). Forniceal DBS, however, bilaterally enhanced dentate neurogenesis in both RTT (DBS-treated RTT: BrdU, $2,012 \pm 269.24$; DCX, $3,146 \pm 340.04$; BrdU/DCX, $1,453 \pm 187.69$) and wild-type mice (DBS-treated wild-type: BrdU, $1,833 \pm 274.79$; DCX, $3,686 \pm 426.92$; BrdU/DCX, $1,526 \pm 257.46$) (Fig. 4e–g and Extended Data Fig. 7b) such that levels of neurogenesis in DBS-treated RTT mice were even higher than those in sham-treated wild-type controls.

Data suggest that decreased cholinergic signalling plays a role in RTT²⁴ and that forniceal stimulation enhances hippocampal memory through cholinergic modulation in rodents³. Therefore, we examined the effect of hippocampal infusion of muscarinic acetylcholine receptor antagonist atropine on the DBS benefit. There was no difference of fear memory between atropine and vehicle-treated groups in either RTT or wild-type mice, suggesting that DBS must benefit memory via additional mechanisms (Extended Data Fig. 8).

DBS has been used to treat both motor and neuropsychiatric disorders in children²⁵. There are a few case reports showing that pallidal DBS can resolve self-injurious behaviour in Lesch-Nyhan syndrome^{26,27}, and hypothalamic DBS reduces aggressive behaviours in patients with intellectual disability²⁸. To our knowledge, however, DBS has never been examined for cognitive benefits in the context of a childhood intellectual disability. In this study, we demonstrate that DBS improves contextual memory retrieval in a mouse model of Rett syndrome. Using a stimulation protocol that mimics clinical

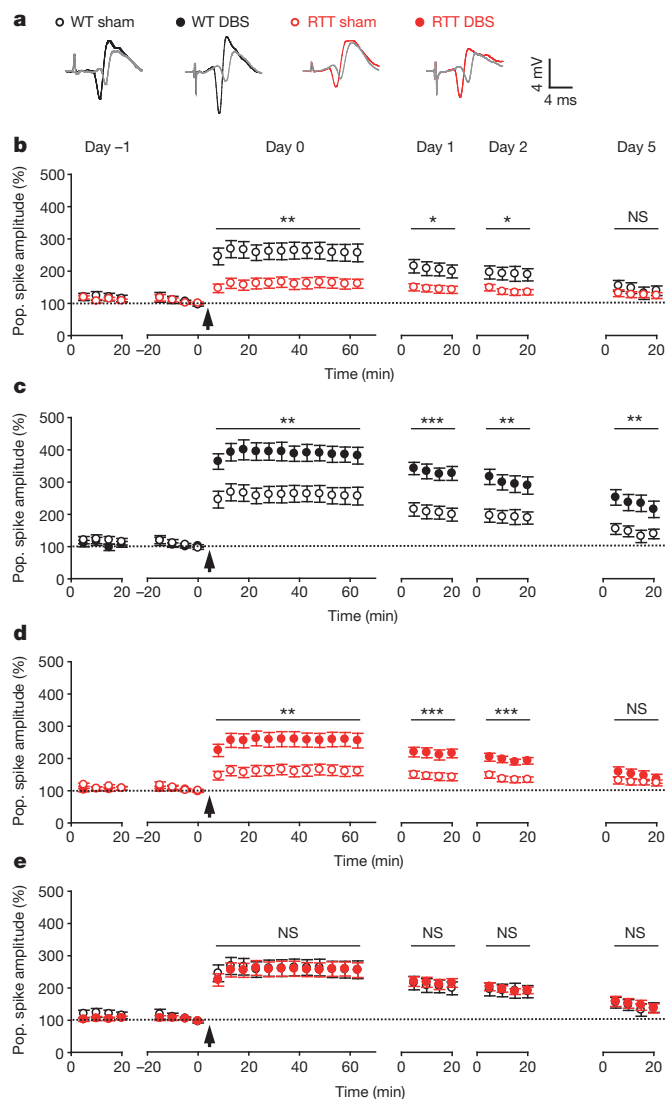


Figure 3 | Forniceal DBS rescues hippocampal synaptic plasticity in freely moving RTT mice. **a**, Superimposed traces of the perforant path recorded in the dentate gyrus 5 min before (grey) and 55 min after (black or red) tetani. **b**, Sham-treated RTT mice ($n = 12$) showed impaired LTP compared to the sham-treated wild-type group ($n = 11$) on day 0 (two-way repeated-measures ANOVA: genotype, $F_{1,21} = 11.34$, $P = 0.003$; time, $F_{15,315} = 40.51$, $P < 0.001$; genotype \times time interaction, $F_{15,315} = 9.36$, $P < 0.001$), day 1 (genotype, $F_{1,21} = 7.46$, $P = 0.012$; time, $F_{3,63} = 5.15$, $P = 0.003$), and day 2 (genotype, $F_{1,21} = 6.50$, $P = 0.019$). **c**, **d**, Forniceal DBS enhanced LTP in both wild-type and RTT mice (DBS-treated wild-type, $n = 12$; sham-treated wild-type, $n = 11$; DBS-treated RTT, $n = 13$; sham-treated RTT, $n = 12$). Two-way repeated-measures ANOVA revealed significant main effects of population spike amplitudes among the four groups on day 0 (group, $F_{3,44} = 17.25$, $P < 0.001$; time, $F_{15,660} = 167.28$, $P < 0.001$; group \times time interaction, $F_{45,660} = 14.50$, $P < 0.001$), day 1 (group, $F_{3,44} = 21.53$, $P < 0.001$; time, $F_{3,132} = 7.69$, $P < 0.001$), day 2 (group, $F_{3,44} = 16.21$, $P < 0.001$; time, $F_{3,132} = 8.96$, $P < 0.001$), and day 5 (group, $F_{3,44} = 8.42$, $P < 0.001$; time, $F_{3,132} = 8.35$, $P < 0.001$). **c**, Forniceal DBS enhanced LTP in wild-type controls on day 0 (treatment, $F_{1,21} = 12.16$, $P = 0.002$; time, $F_{15,315} = 121.93$, $P < 0.001$; treatment \times time interaction, $F_{15,315} = 10.91$, $P < 0.001$), day 1 (treatment, $F_{1,21} = 18.99$, $P < 0.001$; time, $F_{3,63} = 4.77$, $P = 0.005$), day 2 (treatment, $F_{1,21} = 11.25$, $P = 0.003$; time, $F_{3,63} = 3.72$, $P = 0.016$) and day 5 (treatment, $F_{1,21} = 9.44$, $P = 0.006$; time, $F_{3,63} = 6.73$, $P < 0.001$). **d**, Forniceal DBS enhanced LTP in RTT mice on day 0 (treatment, $F_{1,23} = 11.86$, $P = 0.002$; time, $F_{15,345} = 45.02$, $P < 0.001$; treatment \times time interaction, $F_{15,345} = 10.31$, $P < 0.001$), day 1 (treatment, $F_{1,23} = 14.60$, $P < 0.001$; time, $F_{3,69} = 2.91$, $P = 0.041$) and day 2 (treatment, $F_{1,23} = 19.45$, $P < 0.001$; time, $F_{3,69} = 6.09$, $P < 0.001$). **e**, There was no difference of LTP between DBS-treated RTT and sham-treated wild-type mice ($P > 0.05$ for all the test days). Arrows, LTP induction. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$; NS, not significant. Data presented as mean \pm s.e.m.

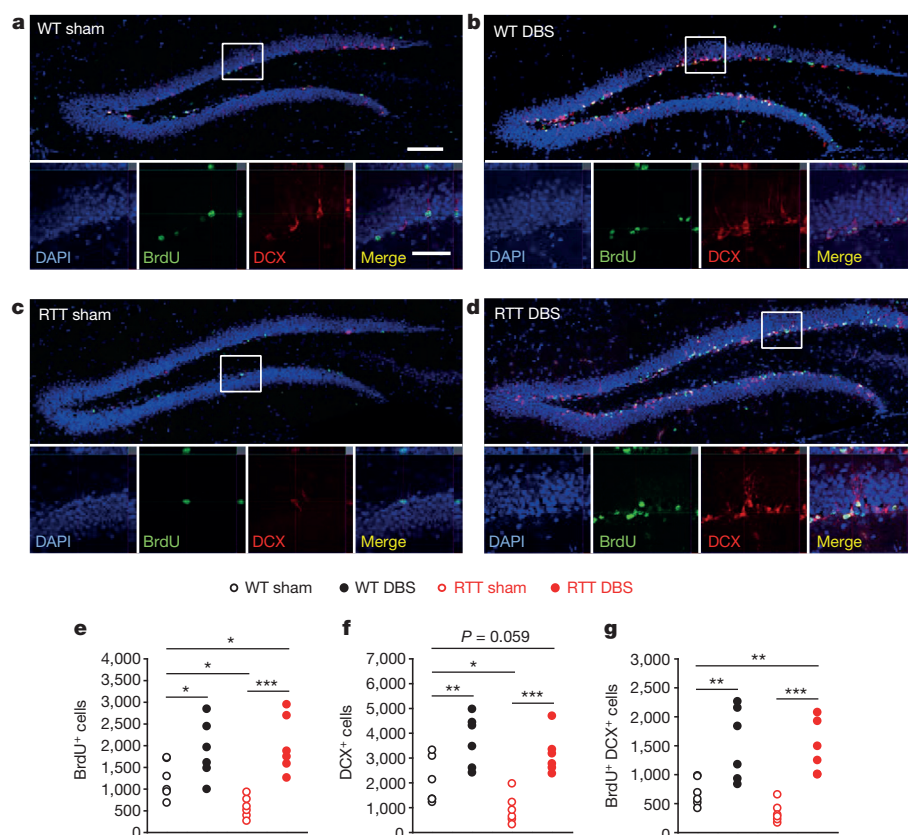


Figure 4 | Forniceal DBS stimulates hippocampal neurogenesis in wild-type and RTT mice. One day after completing the 2 weeks of fornical DBS/sham treatment, animals were perfused for immunohistochemical detection of BrdU- and DCX-positive cells in the dentate gyrus. **a–d**, Representative images ipsilateral to the DBS/sham treatment at low (top; scale bar, 100 μ m) and high (bottom; scale bar, 50 μ m) magnification, showing BrdU⁺ cells (green), DCX⁺ cells (red) and the merge (yellow) from each of the four groups. New neurons were located in the innermost layer of the dentate gyrus. **e–g**, Summary of immunoreactive cell counting ($n = 6$ mice per group). Two-way ANOVA revealed a significant main effect on the numbers of BrdU⁺ cells (**e**; treatment, $F_{1,20} = 23.49$, $P < 0.001$), DCX⁺ cells (**f**; genotype, $F_{1,20} = 5.65$, $P = 0.028$; treatment, $F_{1,20} = 29.65$, $P < 0.001$), and BrdU⁺ DCX⁺ double-staining cells (**g**; treatment, $F_{1,20} = 32.99$, $P < 0.001$). Tukey's post hoc test indicated that sham-treated RTT mice had fewer BrdU⁺ (**e**) and DCX⁺ (**f**) cells than sham-treated wild-type controls. Forniceal DBS increased the numbers of BrdU⁺ (**e**), DCX⁺ (**f**) and BrdU⁺ DCX⁺ double-staining (**g**) cells in DBS-treated wild-type and RTT mice compared to their respective sham-treated controls. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$. Scattergrams show individual values.

treatment, forniceal DBS clearly enhanced contextual fear memory as well as spatial learning and memory in both wild-type and RTT mice. DBS in the fimbria-fornix is so effective in the RTT mice that it restores hippocampus-dependent memory in both tasks to wild-type levels. We also found that forniceal DBS increases hippocampal synaptic plasticity and hippocampal neurogenesis, both of which are central to hippocampal learning and memory^{23,29}.

Although this study was limited to hippocampus-based learning and memory, it is remarkable that DBS could exert any benefit in the face of such profound cognitive impairments as caused by RTT. Future work will explore additional DBS targets to determine the possible benefits of DBS on other RTT features such as dystonia and motor incoordination. Our studies lead us to suggest that DBS should be explored in other models of intellectual disabilities and eventually in human patients, particularly those conditions that cause more focal deficits in learning and memory. Intellectual disabilities as a group affect 2–3% of the population³⁰ and are at present untreatable; their molecular heterogeneity poses a daunting challenge to those looking for viable therapies. The fact that DBS is able to be modulated, reversible, and safe makes it an appealing candidate treatment that could potentially relieve a great deal of suffering.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 8 September 2014; accepted 8 September 2015.

1. Laxton, A. W. *et al.* A phase I trial of deep brain stimulation of memory circuits in Alzheimer's disease. *Ann. Neurol.* **68**, 521–534 (2010).
2. Hamani, C. *et al.* Memory enhancement induced by hypothalamic/fornix deep brain stimulation. *Ann. Neurol.* **63**, 119–123 (2008).
3. Shirvinkar, P. R., Rapp, P. R. & Shapiro, M. L. Bidirectional changes to hippocampal theta-gamma comodulation predict memory for recent spatial episodes. *Proc. Natl Acad. Sci. USA* **107**, 7054–7059 (2010).
4. Phillips, R. G. & LeDoux, J. E. Lesions of the fornix but not the entorhinal or perirhinal cortex interfere with contextual fear conditioning. *J. Neurosci.* **15**, 5308–5315 (1995).

5. Maren, S. & Fanselow, M. S. Electrolytic lesions of the fimbria/fornix, dorsal hippocampus, or entorhinal cortex produce anterograde deficits in contextual fear conditioning in rats. *Neurobiol. Learn. Mem.* **67**, 142–149 (1997).
6. Amir, R. E. *et al.* Rett syndrome is caused by mutations in X-linked *MECP2*, encoding methyl-CpG-binding protein 2. *Nature Genet.* **23**, 185–188 (1999).
7. Chahrour, M. & Zoghbi, H. Y. The story of Rett syndrome: from clinic to neurobiology. *Neuron* **56**, 422–437 (2007).
8. Guy, J., Gan, J., Selfridge, J., Cobb, S. & Bird, A. Reversal of neurological defects in a mouse model of Rett syndrome. *Science* **315**, 1143–1147 (2007).
9. Chao, H. T. *et al.* Dysfunction in GABA signalling mediates autism-like stereotypies and Rett syndrome phenotypes. *Nature* **468**, 263–269 (2010).
10. Moretti, P. *et al.* Learning and memory and synaptic plasticity are impaired in a mouse model of Rett syndrome. *J. Neurosci.* **26**, 319–327 (2006).
11. Samaco, R. C. *et al.* Female *Mecp2*^{+/−} mice display robust behavioral deficits on two different genetic backgrounds providing a framework for pre-clinical studies. *Hum. Mol. Genet.* **22**, 96–109 (2013).
12. Freund, H. J. *et al.* Cognitive functions in a patient with Parkinson-dementia syndrome undergoing deep brain stimulation. *Arch. Neurol.* **66**, 781–785 (2009).
13. Whittle, N. *et al.* Deep brain stimulation, histone deacetylase inhibitors and glutamatergic drugs rescue resistance to fear extinction in a genetic mouse model. *Neuropharmacology* **64**, 414–423 (2013).
14. Suthana, N. *et al.* Memory enhancement and deep-brain stimulation of the entorhinal area. *N. Engl. J. Med.* **366**, 502–510 (2012).
15. Gary-Bobo, E. & Bonvallet, M. Commissural projection to the amygdala through the fimbria fornix system in the cat. *Exp. Brain Res.* **27**, 61–70 (1977).
16. Morris, R. G., Garrud, P., Rawlins, J. N. & O'Keefe, J. Place navigation impaired in rats with hippocampal lesions. *Nature* **297**, 681–683 (1982).
17. Jones, M. W. *et al.* A requirement for the immediate early gene *Zif268* in the expression of late LTP and long-term memories. *Nature Neurosci.* **4**, 289–296 (2001).
18. van Praag, H., Kempermann, G. & Gage, F. H. Running increases cell proliferation and neurogenesis in the adult mouse dentate gyrus. *Nature Neurosci.* **2**, 266–270 (1999).
19. Shors, T. J. *et al.* Neurogenesis in the adult is involved in the formation of trace memories. *Nature* **410**, 372–376 (2001).
20. Stuchlik, A. Dynamic learning and memory, synaptic plasticity and neurogenesis: an update. *Front. Behav. Neurosci.* **8**, 106 (2014).
21. Toda, H., Hamani, C., Fawcett, A. P., Hutchison, W. D. & Lozano, A. M. The regulation of adult rodent hippocampal neurogenesis by deep brain stimulation. *J. Neurosurg.* **108**, 132–138 (2008).
22. Encinas, J. M., Hamani, C., Lozano, A. M. & Enikolopov, G. Neurogenic hippocampal targets of deep brain stimulation. *J. Comp. Neurol.* **519**, 6–20 (2011).
23. Stone, S. S. *et al.* Stimulation of entorhinal cortex promotes adult neurogenesis and facilitates spatial memory. *J. Neurosci.* **31**, 13469–13484 (2011).

24. Wenk, G. L., Naidu, S., Casanova, M. F., Kitt, C. A. & Moser, H. Altered neurochemical markers in Rett's syndrome. *Neurology* **41**, 1753–1756 (1991).
25. Wichmann, T. & DeLong, M. R. Deep brain stimulation for neurologic and neuropsychiatric disorders. *Neuron* **52**, 197–204 (2006).
26. Cif, L. *et al.* Antero-ventral internal pallidum stimulation improves behavioral disorders in Lesch-Nyhan disease. *Mov. Disord.* **22**, 2126–2129 (2007).
27. Deon, L. L., Kalichman, M. A., Booth, C. L., Slavin, K. V. & Gaebler-Spira, D. J. Pallidal deep-brain stimulation associated with complete remission of self-injurious behaviors in a patient with Lesch-Nyhan syndrome: a case report. *J. Child Neurol.* **27**, 117–120 (2012).
28. Franzini, A., Broggi, G., Cordella, R., Dones, I. & Messina, G. Deep-brain stimulation for aggressive and disruptive behavior. *World Neurosurg.* **80**, S29.e11–S29.e14 (2013).
29. Malenka, R. C. & Bear, M. F. LTP and LTD: an embarrassment of riches. *Neuron* **44**, 5–21 (2004).
30. Ramocki, M. B. & Zoghbi, H. Y. Failure of neuronal homeostasis results in common neuropsychiatric phenotypes. *Nature* **455**, 912–918 (2008).

Acknowledgements We thank M. Xue, M. C. Weston and V. Brandt for comments on the manuscript, members of the Zoghbi laboratory for helpful discussions, and

C. M. Spencer, C. T. Wotjak, F. Wei and D. Yu for technical suggestions. This work was supported by the W. M. Keck Foundation (H.Y.Z. and J.T.), the Cockrell Family Foundation, the Rett Syndrome Research Trust, Carl C. Anderson, Sr. and Marie Jo Anderson Charitable Foundation, R01NS057819 (H.Y.Z.), and the Howard Hughes Medical Institute (H.Y.Z.), DP5OD009134 (R.C.S.), R25 N070694 (A.J.P.) and in part by the Neuroconnectivity Core, Mouse Neurobehavioral Core, and Neurovisualization Core of IDDRC at Baylor College of Medicine (U54 HD083092 from the Eunice Kennedy Shriver National Institute of Child Health & Human Development), and the C06RR029965 grant from the National Center for Research Resources.

Author Contributions J.T. and H.Y.Z. designed the experiments. S.H., B.T., Z.W., Y.S., H.T., Y.G., K.U. and J.T. performed the research. S.H., B.T., K.U., H.Y.Z. and J.T. analysed and interpreted the data. R.C.S., A.J.P. and D.J.C. provided comments on the manuscript. S.H., H.Y.Z. and J.T. wrote and edited the paper.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to H.Y.Z. (hzoghbi@bcm.edu) or J.T. (jtang1@bcm.edu).

METHODS

Animals. Adult female *Meep2^{+/-}* mice (14–16 weeks of age at the time of fear conditioning, water maze or LTP test) (Extended Data Fig. 1) on an FVB.129 background were maintained on a 12 h light:12 h dark cycle (light on at 7:00) with standard mouse chow and water *ad libitum* in our on-site AAALAS-accredited facility. They were group-housed up to five mice per cage before surgery and individually housed with nesting material in the cage after surgery in a room maintained at 22 °C. All the experimental procedures and tests were conducted during the light cycle. Behavioural, electrophysiological, immunohistochemical, and pharmacological characterizations of the mice were performed and analysed blind to genotypes and/or treatments. All research and animal care procedures were approved by the Baylor College of Medicine Institutional Animal Care and Use Committee.

Surgery and DBS. Mice were secured on a stereotaxic frame (David Kopf) under 1–2% isoflurane anaesthesia. Bipolar DBS electrodes were constructed with Teflon-coated tungsten wire (bare diameter 50 µm, A-M Systems) and the two tips were horizontally separated by 0.1–0.15 mm. The electrodes were targeted unilaterally to the fimbria–fornix (0.2 mm posterior, 0.2 mm lateral, and 2.3–2.4 mm below the bregma) under the guidance of evoked potentials recorded in the ipsilateral dentate gyrus (1.8–2.0 mm posterior, 1.4–1.6 mm lateral of bregma, 2.2–2.3 mm below the skull³¹). All the electrodes together with the attached connector sockets were fixed on the skull by dental cement. Animals were given at least 2 weeks to recover.

After recovery, mice were assigned into four groups (randomly to DBS or sham groups within the same genotype): wild-type, sham; wild-type, DBS; RTT, sham; and RTT, DBS. Animals in both DBS groups received 1 h DBS daily for 14 consecutive days. Based on DBS parameters widely used in the clinic¹² and the cognitive assessment of DBS in rodents¹³, the DBS was biphasic rectangular pulses (130 Hz, 60 µs pulse duration). This DBS pattern is used both in human subjects^{2,12,32} and rodents¹³. The stimulus intensities were individually optimized to 80% of the threshold that elicits an afterdischarge in the hippocampus^{3,14}. Animals in the sham groups experienced the same procedures as those in the DBS groups except without DBS. There was a 3-week interval between the last DBS administration and the beginning of fear conditioning, water maze, or *in vivo* LTP tests²³. All the mice received DBS/sham treatment and those used for fear conditioning and EEG recordings were habituated to the headstage/wiring system in an environment different from both conditioning chamber and the cue memory test cage for 3 days (20 min per day) immediately before these tests.

After finishing all experiments, mice were euthanized with an overdose of isoflurane. An anodal current (30 µA, 10 s) was passed through the electrode wire to verify the electrode placements. Frozen 30-µm coronal sections were cut and stained with cresyl violet.

Behavioural tests. Tests of light/dark chamber and open field were conducted 1 week before fear conditioning or water maze. Wire-hang, dowel-walk, rotarod and three-chamber tests, as well as a test of pain threshold, were performed 1–2 weeks after fear conditioning or water maze. (Extended Data Fig. 1) For each test, mice were given at least 30 min to habituate after transport to the behavioural testing room before any tests were conducted. The light intensity of 150 lx and the background white noise at 60 dB were presented during the habituation and throughout the testing periods.

Fear conditioning. A delayed fear conditioning protocol was employed to evaluate hippocampus-dependent contextual fear memory and hippocampus-independent cue fear memory. On day 0 animals were trained in a mouse fear conditioning chamber with a grid floor that could deliver an electric shock (Med Associates, Inc.). This enclosure was located in a sound-attenuating box that contained a digital camera, a loudspeaker and a house light. Each mouse was initially placed in the chamber and left undisturbed for 2 min, after which a tone (30 s, 5 kHz, 80 dB) coincided with a scrambled foot shock (2 s, 0.7 mA). The tone/foot-shock stimuli were repeated after 1 min. After an additional 1 min, the mouse was removed and returned to its home cage. Fear memory retention was evaluated 3 h, 1 day, 3 days, and 7 days after training unless stated otherwise. At each time point, mice were first recorded for 3 min in the same chamber (cleaned with 70% ethanol) without tone. The mice were then tested in a novel cage (cleaned with 1% acetic acid) where a 3-min tone was presented after the animals had acclimated to the cage for 3 min. Mouse behaviour was recorded and scored automatically by ANY-maze (Stoelting). Freezing, defined as an absence of all movement except for respiration³³, was scored only if the animal was immobile for at least 1 s³⁴. The percentage of time spent freezing during the tests serves as an index of fear memory. Cued fear memory was the subtraction of freezing time between the tone phase and the no-tone phase. Data are shown as mean ± s.e.m. and analysed by two-way ANOVA followed by Tukey's post hoc analysis.

Morris water maze. The Morris water maze was used to assess spatial learning and memory in RTT mice and the effect of fornical DBS. This assay was performed as previously described with a few modifications^{35,36}. The pool (120 cm in diameter) was filled with water (50 cm deep, 22–24 °C) made opaque with non-toxic white tempera paint. Visual cues were set on the wall of the testing room, at least 1 m from the pool edge. The ANY-maze tracking system (Stoelting) was used to track and analyse animal swimming. Mice were tested for their ability to find an escape platform (10 cm in diameter) on three different components of the test: visible platform acquisition, hidden platform acquisition, and subsequent probe trial in the absence of the platform. In each case, the criterion for learning is an average latency of 15 s or less to locate the platform across a block of four consecutive trials (15 s interval) per day. Mice are given up to 9 days to reach this criterion for learning. Twenty-four hours after the last training trial the mice were given a probe trial. During the probe trial, the platform was removed, and each animal was allowed 60 s to search the pool. The amount of time and distance that each animal spent in each quadrant was recorded. The number of times a subject crossed the imaging location of the platform during training (platform crossing) was determined. Data of memory acquisition are expressed as mean ± s.e.m. and analysed by two-way ANOVA with repeated-measures followed by Tukey's post hoc analysis. Data of probe tests are shown as individual values and two-tailed *t*-test is used to compare the searching time in the target quadrants and the platform crossing numbers between groups, and two-tailed paired *t*-test for comparing the searching time in quadrants within groups.

Open field. The open field apparatus consists of a clear, open Plexiglas box (40 × 40 × 30 cm, Stoelting) with overhead camera and photo beams to record horizontal and vertical movements. Activity was quantified over a 30-min period by ANY-maze (Stoelting). Data are shown as mean ± s.e.m. and analysed by two-way ANOVA followed by Tukey's post hoc analysis.

Light–dark box. The light–dark box assay was performed as published with few modifications³⁷. The box consisted of a clear Plexiglas chamber (40 × 20 × 30 cm) with an open top separated from a covered black chamber (40 × 20 × 30 cm) by a black partition with a small opening (Stoelting). RTT and wild-type mice were placed into the illuminated side and allowed to explore freely for 10 min. An ANY-maze system with photo beam and overhead camera (Stoelting) was used to score the mice for the number and latency of entries and the time spent in each compartment. The mouse must place 50% of body length into either the light or dark compartment to be scored as an entry. Data are shown as mean ± s.e.m. and analysed by two-way ANOVA with Tukey's post hoc analysis.

Wire hang and dowel walk. These assays were performed as previously published with a few modifications³⁸. RTT and wild-type mice were tested for motor coordination. For the wire-hang test the mouse was held by the tail and allowed to grasp with its forepaws the middle of a single 3-mm plastic-coated wire suspended 15 inches above a plastic-covered foam pad and released. For the dowel test the mouse was placed onto a 0.5-inch wooden dowel suspended 15 inches above a plastic-covered foam pad, and the total number of side touches and latency to fall were measured with a 120 s cut-off. Data are shown as mean ± s.e.m. and analysed using two-way ANOVA with Tukey's post hoc analysis.

Accelerating rotarod. This assay was performed as previously published with a few modifications³⁸. RTT and wild-type mice were placed on the rotating cylinder of an accelerating rotarod apparatus (Ugo Basile) and allowed to move freely as the rotation increased from 5 r.p.m. to 40 r.p.m. over a 5-min period. Latency to fall was recorded when the mouse fell from the rod or when the mouse had ridden the rotating rod for two revolutions without regaining control. Data are shown as mean ± s.e.m. Latency to fall was analysed by two-way repeated-measures ANOVA with Tukey's post hoc analysis.

Three-chamber interaction. The test was performed as previously described^{35,39}. RTT and wild-type mice were used in this assay. For the habituation stage, test mice were placed in the middle chamber of the three-chamber apparatus (Ugo Basile) equipped with two empty, barred cages in the corners of the left and right chambers. They were allowed to explore freely for 10 min, with their movement tracked and recorded using the ANY-maze software pack (Stoelting), and interaction time with each cage scored by an investigator blind to genotype and treatment group. For the social versus object stage, an age- and size-matched C57BL/6 female mouse was placed in one cage and a Lego block of similar size was placed in the other cage. The test mouse was again placed in the middle chamber and allowed to explore freely for 10 min, with movement and interaction time recorded as before. Interaction time and time in each zone are shown as mean ± s.e.m. and analysed by two-way ANOVA with Tukey's post hoc analysis.

Pain threshold. The test was performed as previously published with a few modifications⁴⁰. At the end of the test battery, animals were placed into the conditioning chamber. Every 30 s, a 2-s scrambled electric foot shock with 0.05 mA increments (starting from 0 mA) was applied. The shock current thresholds of

flinch, vocalization, and jumping were each recorded. Data are shown as mean \pm s.e.m. and analysed by two-way ANOVA with Tukey's post hoc analysis.

Induction and recording of hippocampal synaptic plasticity *in vivo*. To determine the effect of fornical DBS on hippocampal synaptic plasticity, an additional concentric stimulating electrode was implanted ipsilaterally in the medial perforant path (0.2 mm posterior and 2.8–3.0 mm lateral of lambda, 1.0–1.3 mm below the dura). The stimulating and recording electrodes were surgically implanted as previously described^{41,42} with the following modifications. The final depth of the electrodes was determined by electrophysiological guidance. A cortical silver ball, placed contralaterally, served as a recording reference as well as ground. Dental cement was used to anchor the electrode assembly that is connected to a unity gain preamplifier, and the connecting device for chronic recordings. After recovery from surgical implantation, mice were transported and habituated to the recording system during each of the 4 days before starting the LTP test. Signals were amplified (100 \times), filtered (bandpass, 0.1–5 kHz), digitized at 10 kHz, and stored on disk for off-line analysis (pClamp10 and 1440A; Molecular Devices). To evaluate whether fornical DBS influence the input-output (I/O) relation in the perforant path–dentate pathway, I/O curves were generated for each mouse 1 day before and 3 weeks after the DBS/sham treatment. For LTP evaluation, test responses elicited by 0.033 Hz monophasic pulses (0.1 ms duration) were recorded for 20-min periods on consecutive days at an intensity that evoked 40% of the maximal population spike. Following 2 days of stable baseline, a tetanus was delivered to the perforant path for LTP induction. Pulse width was doubled during tetani, which consisted of six series of six trains of six stimuli at 400 Hz, 200 ms between trains, 20 s between series. Responses were measured for 60 min after tetanus and again for 20 min at 24 h, 48 h and 120 h after tetanus. Since the latency of the population spike usually decreases following LTP induction, it is impractical to compare the initial slope of the fEPSP (field excitatory postsynaptic potential) before and after LTP induction in awake animals^{17,43}. Accordingly, we quantified the amplitude of the population spikes⁴². Data were averaged every 5 min and normalized to the baseline measured over the 10 min before tetanic stimulation and presented as mean \pm s.e.m. Two-way repeated-measures ANOVA was used for data analysis.

Recordings of hippocampal local field potentials (LFPs) and data analysis. The recording electrode of LFPs was targeted to the upper molecular layer of the dentate gyrus with the reference electrode in the corpus callosum. Recording of LFPs was conducted under matched behavioural states for RTT and wild-type mice. Signals were amplified (100 \times), filtered (bandpass, 0.1–5 kHz), digitized at 2 kHz, and stored on disk for off-line analysis (pClamp10 and 1440A; Molecular Devices). The power spectrum of the LFPs was calculated at 0.244 Hz resolution, using the built-in function of pClamp 10. Then the relative power of hippocampal theta activity was normalized as the ratio between the power of the theta signal at 4.15–11.96 Hz and the power of the signal at 0–100 Hz^{44,45}. LFP recordings were excluded from further analysis if mice made large movements during recordings or electrical artefacts were evident³.

Intracranial drug infusion and histology. Under isoflurane anaesthesia, custom-made 26G guide cannulas were implanted bilaterally (1.8 mm posterior, 1.2 mm lateral, and 1.2 mm below the bregma) in RTT mice and littermate wild-type controls. Other than the DBS electrode, a recording electrode (Teflon-coated tungsten wire, 50 μ m bare diameter) was unilaterally attached to the guide cannula for the recordings of evoked potential during surgical implantation and LFPs in the dentate. After 2 weeks of recovery, mice were randomly assigned into drug- or vehicle-treated groups. Atropine sulphate (Sigma) was dissolved in PBS (pH 7.4) and the solution was back-filled into a 33G injector. The solution of 0.5 μ l drug (1.0 μ g atropine)^{46,47} or PBS vehicle per side was microinfused over 1 min into the dorsal hippocampus (dHP) through a pump (Harvard Apparatus)⁴². The tip of the injector was 1.0 mm below the guide tip and \sim 0.6 mm from the tip of the recording electrode. Following each injection, the injector was left in place for an additional minute to allow drug diffusion. To determine the effect of atropine on hippocampal potentials evoked by FFX stimulation, atropine was infused into the dHP after 20 min of baseline test (0.033 Hz) and the responses were followed up for another 30 min in freely moving mice. To determine the effect of atropine on hippocampal theta activity, dentate LFPs were recorded for 5 min immediately before and another 5 min after atropine/vehicle infusion (10–15 min after the infusion). To evaluate the influence of atropine on DBS effect, RTT and wild-type mice receiving fornical DBS were infused with atropine or vehicle 15 \pm 2 min before fear conditioning training. Then fear memory was tested 24 h after training.

At the end of the experiment, 4% methylene blue in PBS (0.2 μ l) was injected to each injection site. Mice were euthanized with an overdose of isoflurane. An anodal current (30 μ A, 10 s) was passed through the electrode wire to verify electrode placements. Frozen 30- μ m coronal sections were cut and stained with haematoxylin. Mice with blocked guide cannula(s) or with injection site(s) outside of the dorsal hippocampus were excluded from data analysis.

Immunohistochemistry. To check whether fornical DBS increased the activity of dentate neurons, we assessed expression of the *Fos* gene 60 min after 2 h of DBS or sham treatment. To evaluate whether DBS increases neurogenesis in the dentate gyrus, mice were injected with 5-bromo-2'-deoxyuridine (BrdU, Sigma-Aldrich, 75 mg kg⁻¹, i.p.) immediately after daily DBS or sham treatment for 12 days from the third day of the 2 weeks of DBS/sham treatment.

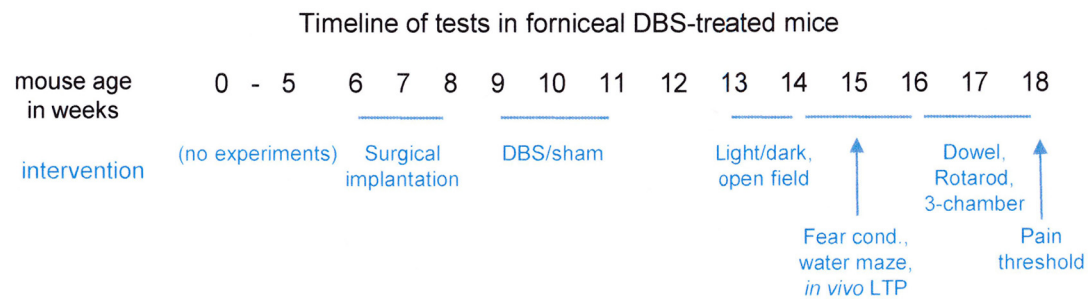
Immunohistochemical staining. Mice were euthanized with an overdose of isoflurane and subsequently perfused with 50 ml 0.1 M phosphate-buffered saline (PBS, pH 7.4) followed by 60 ml cold 4% paraformaldehyde in 0.1 M PBS. Brains were removed, post fixed with the same fixative for 12 h, and then transferred to 30% sucrose at 4 °C. One of every six of the floating serial 50- μ m coronal sections throughout the hippocampus was processed for c-Fos, DCX, and BrdU immunohistochemistry. After three rinses with 0.1 M PBS, the sections were incubated in blocking solution (0.3% Triton X-100, 5% normal goat serum in 0.1 M PBS) for 1 h at room temperature followed by 48 h incubation with primary antibody at 4 °C. Primary antibodies and their final concentrations were as follows: rabbit polyclonal anti-c-Fos (1:10,000, PC38 Millipore), rabbit anti-DCX (1:400, Cell Signaling Technology) and rat monoclonal anti-BrdU (1:140, OBT0030, Accurate Chemical & Scientific Corporation). Sections were incubated with Alex Fluor 568 goat anti-rabbit IgG (1:500, A11036, Invitrogen) and/or Alex Fluor 488 goat anti-rat IgG (1:500, A11006, Invitrogen) for 2 h and counterstained with DAPI (900 nM, D1036, Invitrogen) for 15 min at room temperature in darkness. The sections were washed three times with 0.1 M PBS and mounted using Fluoromount-G (Southern Biotech). For c-Fos staining, sections were pre-treated with 0.3% H₂O₂ in PBS for 30 min at room temperature. For BrdU detection, sections were pre-treated with 2 M HCl for 30 min at 37 °C and washed in 0.1 M borate buffer, pH 8.4, for 10 min.

Imaging and quantification. Z stacks of 2- μ m thick single-plane images at 20 \times magnification were collected through the entire thickness of the slice by employing a laser scanning microscope LSM 710 (Carl Zeiss). Each slice has 12–13 z-axis optical slices; the sixth optical imaging was selected for counting c-Fos-, BrdU- and DCX-positive cells. Digital images were routed into a PC for quantitative analyses using ImageJ software (NIH). Double labelling for BrdU- and DCX-positive cells was assessed through the entire z axis of each cell. For quantification, six sections per mouse brain were counted. Resulting cell numbers were multiplied by 6 to obtain the estimated total number of positive cells per dentate gyrus.

Statistical analyses. Sample sizes of mice were determined based on prior statistics and data phenotyping the RTT mice¹¹. Animals with disconnected electrode implants before the completion of the experiments were excluded from data analyses. Data were analysed using two-way repeated-measures ANOVA. If any of the main effects were significant, Tukey's post hoc analysis was used for all pairwise multiple comparisons unless otherwise specified. In all cases, $P < 0.05$ was set as the cut-off for statistical significance. SigmaPlot 12 was used to create all the summarized plots as well as all the statistical tests in this study.

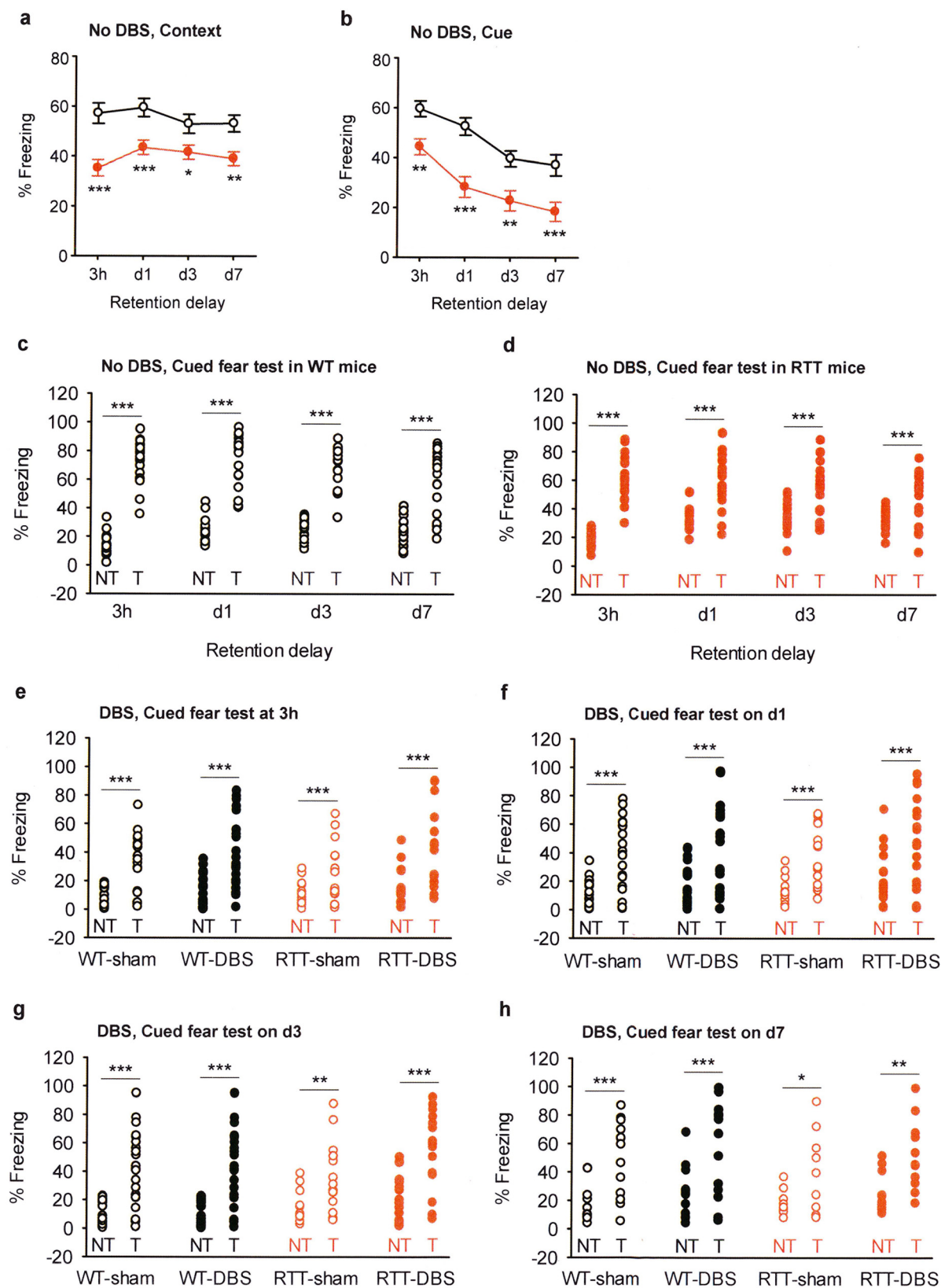
- Paxinos, G. & Franklin, K. B. J. *The Mouse Brain in Stereotaxic Coordinates* (Academic Press, 2001).
- Laxton, A. W., Lipsman, N. & Lozano, A. M. Deep brain stimulation for cognitive disorders. *Handb. Clin. Neurol.* **116**, 307–311 (2013).
- Maren, S., De Oca, B. & Fanselow, M. S. Sex differences in hippocampal long-term potentiation (LTP) and Pavlovian fear conditioning in rats: positive correlation between LTP and contextual learning. *Brain Res.* **661**, 25–34 (1994).
- Corcoran, K. A. & Maren, S. Hippocampal inactivation disrupts contextual retrieval of fear memory after extinction. *J. Neurosci.* **21**, 1720–1726 (2001).
- Moy, S. S. *et al.* Mouse behavioral tasks relevant to autism: phenotypes of 10 inbred strains. *Behav. Brain Res.* **176**, 4–20 (2007).
- Huang, H. S. *et al.* Behavioral deficits in an Angelman syndrome model: effects of genetic background and age. *Behav. Brain Res.* **243**, 79–90 (2013).
- Bouwknegt, J. A. & Paylor, R. Behavioral and physiological mouse assays for anxiety: a survey in nine mouse strains. *Behav. Brain Res.* **136**, 489–501 (2002).
- Shahbazian, M. *et al.* Mice with truncated MeCP2 recapitulate many Rett syndrome features and display hyperacetylation of histone H3. *Neuron* **35**, 243–254 (2002).
- Nadler, J. J. *et al.* Automated apparatus for quantitation of social approach behaviors in mice. *Genes Brain Behav.* **3**, 303–314 (2004).
- Marsicano, G. *et al.* The endogenous cannabinoid system controls extinction of aversive memories. *Nature* **418**, 530–534 (2002).
- Davis, S., Bliss, T. V., Dutrieux, G., Laroche, S. & Errington, M. L. Induction and duration of long-term potentiation in the hippocampus of the freely moving mouse. *J. Neurosci. Methods* **75**, 75–80 (1997).
- Tang, J. & Dani, J. A. Dopamine enables *in vivo* synaptic plasticity associated with the addictive drug nicotine. *Neuron* **63**, 673–682 (2009).
- Malleret, G. *et al.* Inducible and reversible enhancement of learning, memory, and long-term potentiation by genetic inhibition of calcineurin. *Cell* **104**, 675–686 (2001).
- Nokia, M. S., Anderson, M. L. & Shors, T. J. Chemotherapy disrupts learning, neurogenesis and theta activity in the adult brain. *Eur. J. Neurosci.* **36**, 3521–3530 (2012).

45. Nokia, M. S., Sisti, H. M., Choksi, M. R. & Shors, T. J. Learning to learn: theta oscillations predict new learning, which enhances related learning and neurogenesis. *PLoS ONE* **7**, e31375 (2012).
46. Jafari-Sabet, M. NMDA receptor blockers prevents the facilitatory effects of post-training intra-dorsal hippocampal NMDA and physostigmine on memory retention of passive avoidance learning in rats. *Behav. Brain Res.* **169**, 120–127 (2006).
47. Jiao, R., Yang, C., Zhang, Y., Xu, M. & Yang, X. Cholinergic mechanism involved in the nociceptive modulation of dentate gyrus. *Biochem. Biophys. Res. Commun.* **379**, 975–979 (2009).



Extended Data Figure 1 | Timeline of forniceal DBS tests in RTT and wild-type mice.

○ WT ● RTT

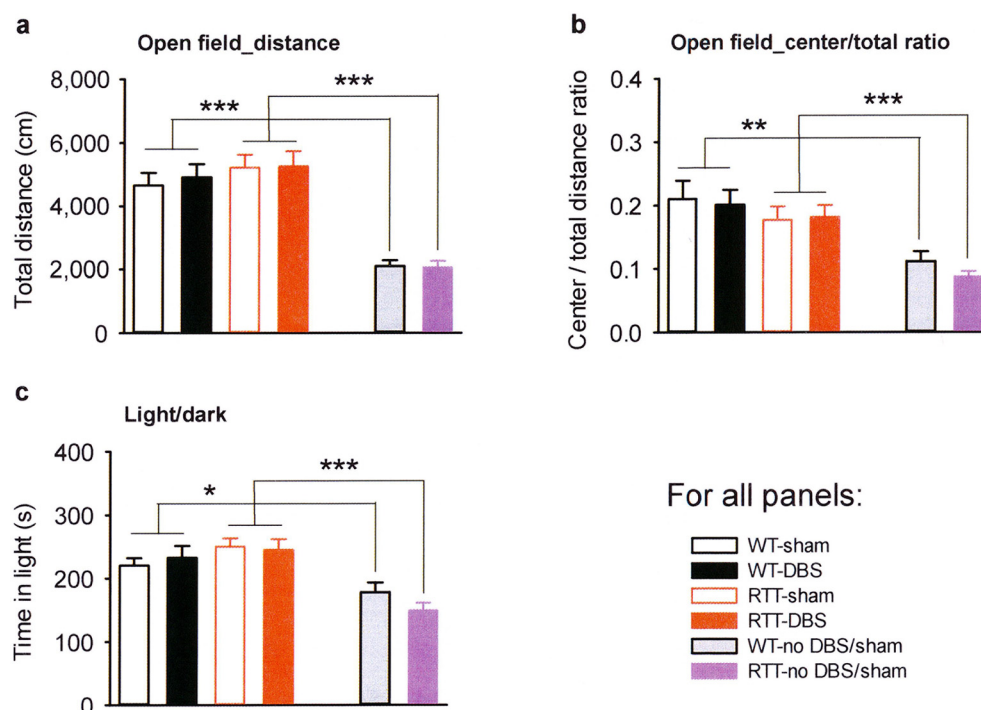


Extended Data Figure 2 | Fear memory in RTT mice and wild-type control animals. All mice were trained with tone–foot-shock pairings on day 0.

Memory retention was tested 3 h, 1 day, 3 day, and 7 day after training.

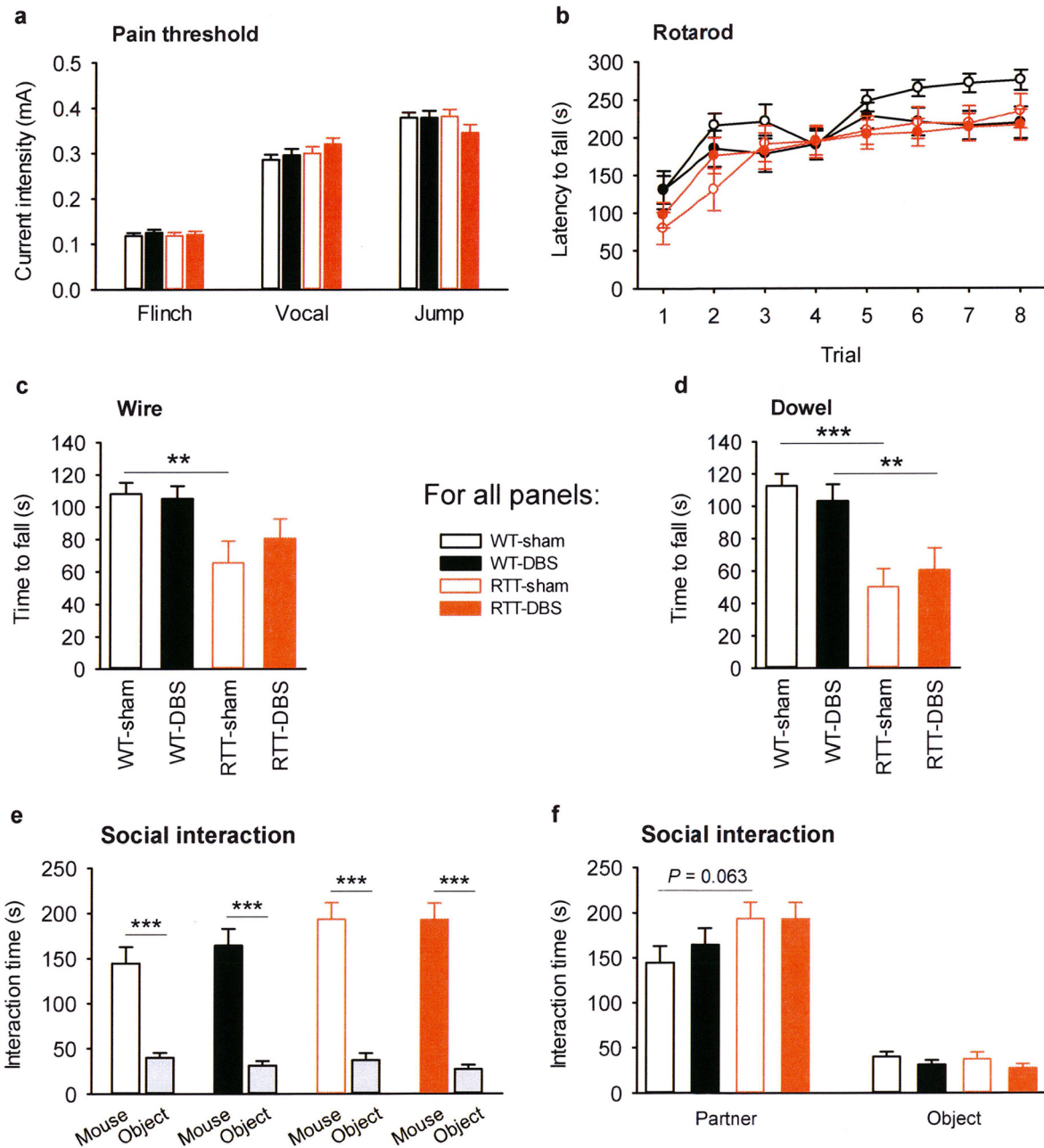
a, b, Impaired fear memory in RTT mice ($n = 20$) compared to wild-type (WT) littermates ($n = 20$). These animals were implanted with electrodes but did not receive DBS or sham treatment. A significant main effect of genotype was observed (two-way repeated-measures ANOVA followed by Tukey's post hoc test: context, $F_{1,38} = 15.32$, $P < 0.001$; cue, $F_{1,38} = 20.70$, $P < 0.001$). * $P < 0.05$; ** $P < 0.01$; *** $P < 0.001$ versus wild type. **c, d,** Cued fear memory in RTT mice ($n = 20$) and wild-type littermates ($n = 20$) that were implanted with

electrodes but without DBS or sham treatment. During the retention test, freezing in the tone phase (T) was significantly more than in the no tone phase (NT) across all the test time points in both wild-type (**c**) and RTT mice (**d**). **e–h,** Retrieval of cue fear memory in DBS- or sham-treated RTT and wild-type mice. During the cued memory test, all four groups of animals actively responded to the tone presentation (WT-sham, $n = 21$; WT-DBS, $n = 21$; RTT-sham, $n = 14$; RTT-DBS, $n = 17$). There was a significant increase of freezing time in the tone phase (T) compared to the no-tone phase (NT) at each of the test time points over all the groups. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$ (two-tailed paired t -test). All data are presented as mean \pm s.e.m.



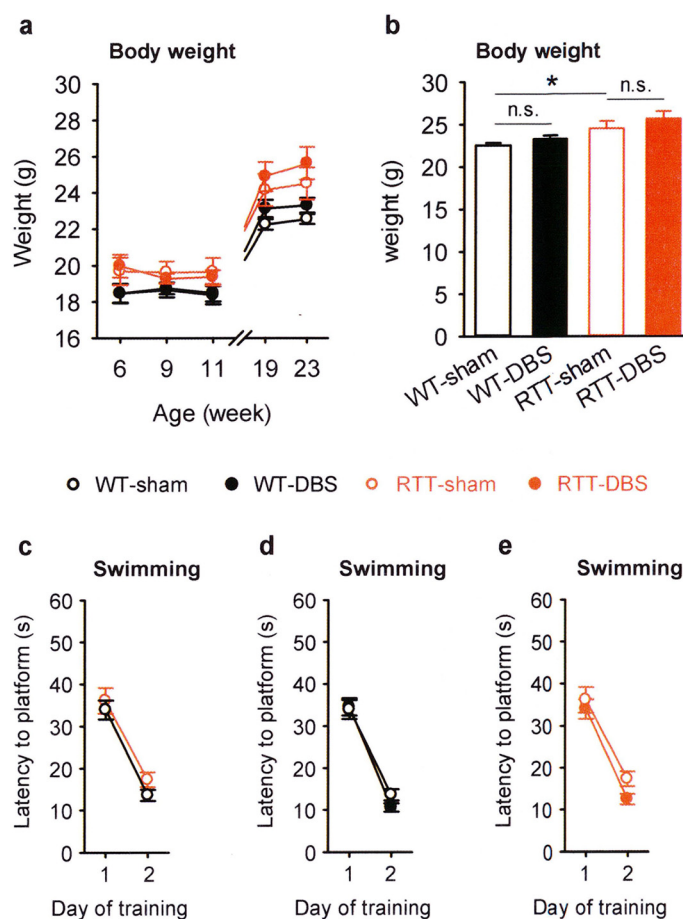
Extended Data Figure 3 | Increased handling, but not forniceal DBS, increased locomotor activity and decreased the anxiety level in RTT and wild-type mice. **a**, There was no difference among the four DBS/sham-treated groups in the total distance travelled in the open-field test (WT-sham, $n = 20$; WT-DBS, $n = 20$; RTT-sham, $n = 17$; RTT-DBS, $n = 18$; genotype, $F_{1,71} = 1.13$, $P = 0.292$; treatment, $F_{1,71} = 0.13$, $P = 0.724$; genotype \times treatment, $F_{1,71} = 0.063$, $P = 0.803$). RTT and wild-type mice that received DBS/sham treatment travelled longer distances than RTT ($n = 20$) and wild-type ($n = 20$) animals that were implanted with electrodes but did not experience DBS/sham procedures, respectively. **b**, During the open-field test, there was no difference in the centre:total distance ratio among the four DBS groups (genotype, $F_{1,71} = 1.22$, $P = 0.273$; treatment, $F_{1,71} = 0.0079$, $P = 0.93$;

genotype \times treatment, $F_{1,71} = 0.081$, $P = 0.777$). Both RTT and wild-type mice that received DBS/sham treatment travelled more in the centre area compared to implanted RTT and wild-type animals that did not receive DBS/sham procedures. **c**, In the light/dark test there was no difference in the amount of time spent in the light compartment among the four chronically treated groups ($n = 12$ per group; two-way ANOVA: genotype, $F_{1,44} = 1.83$, $P = 0.183$; treatment, $F_{1,44} = 0.057$, $P = 0.813$; genotype \times treatment, $F_{1,44} = 0.33$, $P = 0.567$). Both RTT and wild-type mice that received DBS/sham treatment spent more time in the light compartment than implanted RTT ($n = 15$) and wild-type ($n = 14$) animals that did not receive DBS/sham procedures. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$ (two-tailed t -test). All data are presented as mean \pm s.e.m.



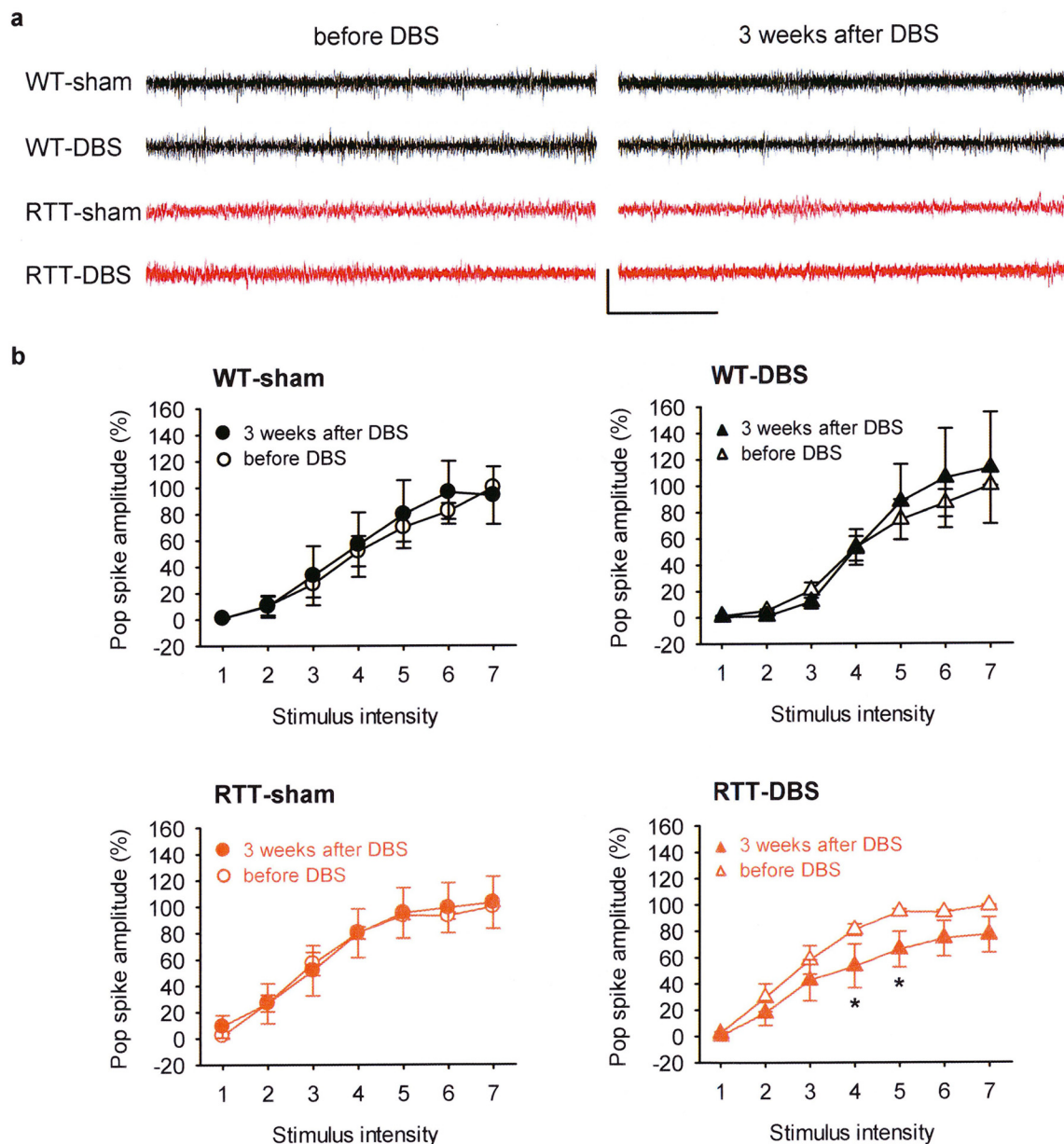
Extended Data Figure 4 | Forniceal DBS did not alter the pain threshold, motor function or social behaviour in RTT or wild-type mice. **a**, There was no group difference in foot shock threshold intensities to evoke flinch, vocalization or jumping (WT-sham, $n = 14$; WT-DBS, $n = 14$; RTT-sham, $n = 11$; RTT-DBS, $n = 12$; two-way ANOVA, no significant main effect of genotype, treatment, or genotype \times treatment interaction, $P > 0.05$). **b**, In a rotarod test ($n = 12$ mice per group), latency to fall increased over trials but there was no difference among the four groups (two-way repeated measures ANOVA: group, $F_{3,44} = 1.68$, $P = 0.184$; trial, $F_{7,308} = 34.26$, $P < 0.001$; group \times trial interaction, $F_{21,308} = 1.22$, $P = 0.230$). **c**, RTT mice showed decreased latency to fall in the wire-hang test compared to wild-type animals, but there was no difference between DBS- and sham-treated groups for either RTT or wild-type mice ($n = 12$ per group; two-way ANOVA: genotype, $F_{1,44} = 10.41$, $P = 0.002$; treatment, $F_{1,44} = 0.33$, $P = 0.566$; genotype \times treatment interaction, $F_{1,44} = 0.75$, $P = 0.392$). **d**, RTT mice

showed a decreased latency to fall in the dowel test compared to wild-type animals, but there was no difference between DBS- and sham-treated groups for either genotype ($n = 12$ per group; genotype, $F_{1,44} = 23.63$, $P < 0.001$; treatment, $F_{1,44} = 0.0018$, $P = 0.966$; genotype \times treatment interaction, $F_{1,44} = 0.83$, $P = 0.367$). **e**, **f**, In the three chamber test, all four groups of animals ($n = 12$ per group) showed a clear preference for the partner mice compared to the object (**e**). Two-way ANOVA revealed a significant genotype main effect of the interaction time with the partner mice ($F_{1,44} = 4.56$, $P = 0.038$), indicating altered social behaviour in RTT mice ($P = 0.063$, RTT-sham versus WT-sham, Tukey's post hoc). However, DBS did not change the interaction time with the partners (treatment, $F_{1,44} = 0.28$, $P = 0.597$; genotype \times treatment interaction, $F_{1,44} = 0.31$, $P = 0.579$) or the object (treatment, $F_{1,44} = 2.64$, $P = 0.111$; genotype \times treatment interaction, $F_{1,44} = 0.015$, $P = 0.905$) (**f**). $**P < 0.01$, $***P < 0.001$ (Tukey's post hoc in **c**, **d**; two-tailed paired t -test in **e**). All data are presented as mean \pm s.e.m.



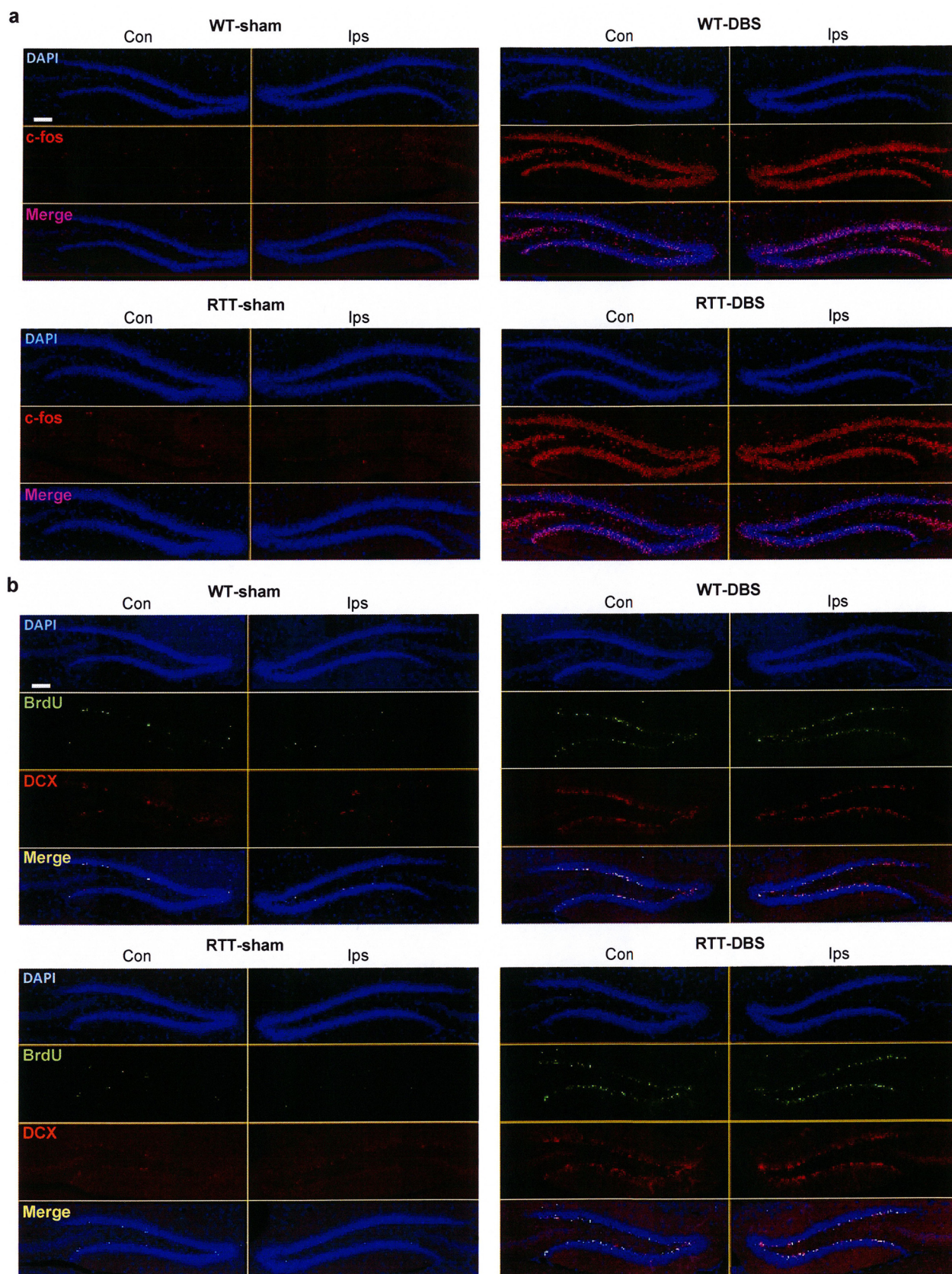
Extended Data Figure 5 | Forniceal DBS did not alter the body weight, visual or sensorimotor skills in RTT or wild-type mice. **a**, All four groups ($n = 12$ mice per group) showed changes in body weight over time. Two-way repeated measure ANOVA revealed a significant main effect of group ($F_{3,44} = 6.73$, $P < 0.001$) and age ($F_{4,176} = 89.32$, $P < 0.001$). Tukey's post hoc showed that sham-treated RTT mice were significantly heavier than sham-treated wild-type mice ($P = 0.015$), but there was no difference in body weight between sham-treated and DBS-treated wild-type mice ($P = 0.861$) or between sham-treated and DBS-treated RTT mice ($P = 0.099$). **b**, Comparison of body weight at the age of 23 weeks among the four groups (two-way ANOVA: genotype, $F_{1,44} = 10.06$, $P = 0.003$; treatment: $F_{1,44} = 1.93$,

$P = 0.172$). **c–e**, Swimming test in the water maze task with a flagged platform ($n = 18$ mice per group). Sham-treated RTT mice did not have different escape latencies than sham-treated wild-type controls (**c**, two-way repeated-measures ANOVA: genotype, $F_{1,34} = 1.73$, $P = 0.197$; genotype \times treatment interaction, $F_{1,34} = 0.133$, $P = 0.718$). DBS did not change the escape latencies in either wild-type controls (**d**; treatment, $F_{1,34} = 0.44$, $P = 0.513$; treatment \times day interaction, $F_{1,34} = 1.24$, $P = 0.273$) or RTT mice (**e**, treatment, $F_{1,34} = 2.36$, $P = 0.134$; treatment \times day interaction, $F_{1,34} = 0.41$, $P = 0.524$). * $P < 0.05$; n.s., not significant (Tukey's post hoc). All data are presented as mean \pm s.e.m.



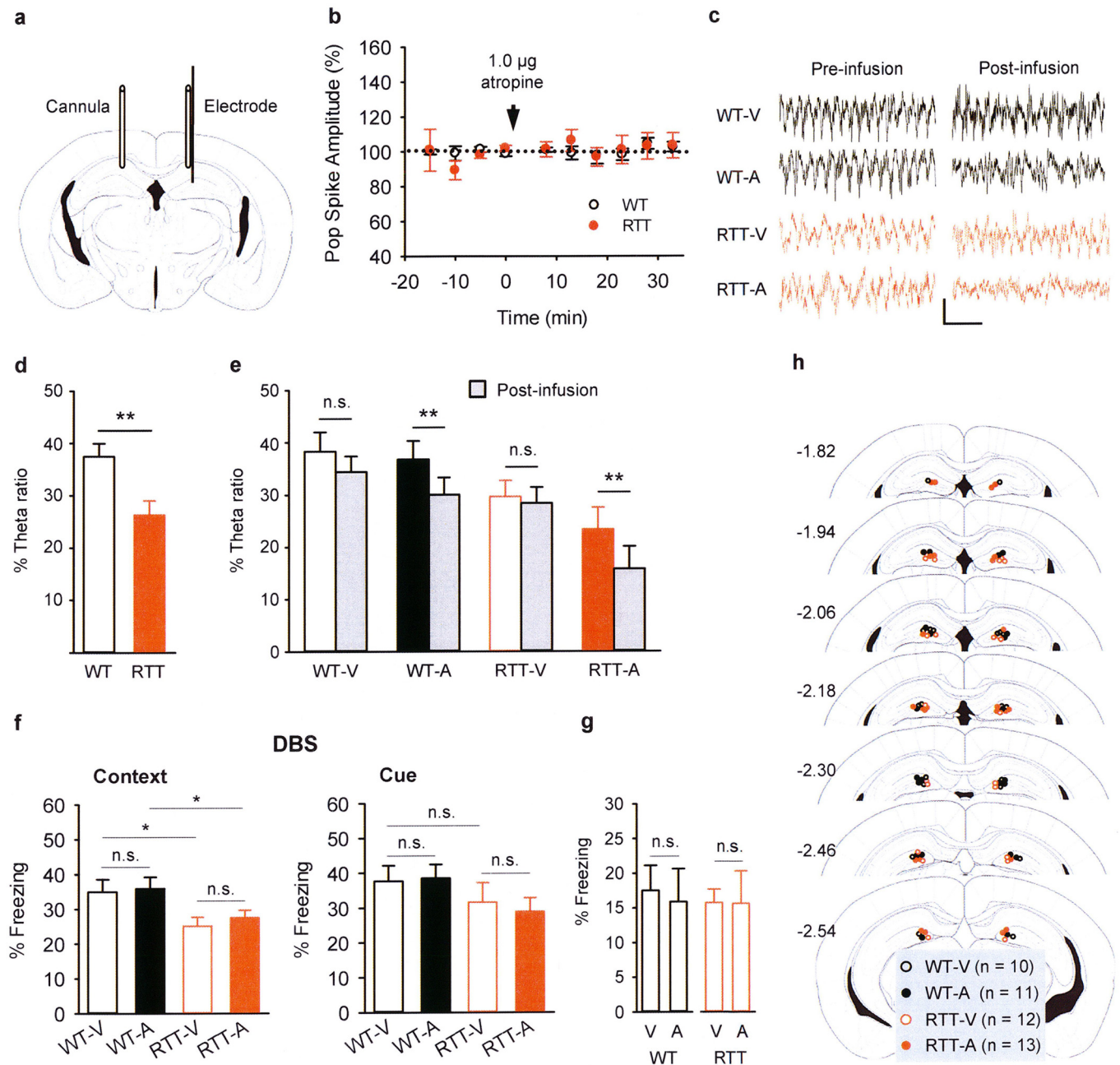
Extended Data Figure 6 | Effect of forniceal DBS on hippocampal electrophysiological signatures. **a**, Representative traces of LFPs recorded in the dentate gyrus 1 day before and 3 weeks after DBS/sham treatment. There were no electrographic seizure spikes in any of the four groups of mice after DBS/sham treatment. Scale bars: 10 s, 1 mV. **b**, Input-output (I/O) curves of the evoked responses of the perforant path recorded in the dentate gyrus in DBS/sham-treated mice. For each of the four groups, I/O curves were generated 1 day before and 3 weeks after forniceal DBS. All data points were normalized to the maximum value of the population spike amplitude before DBS/sham and

the abscissa represents the seven increments used in each mouse. The I/O relationship was not altered by DBS in sham-treated wild-type mice (WT-sham; $n = 5$, $F_{1,4} = 0.062$, $P = 0.818$), DBS-treated wild-type mice (WT-DBS; $n = 4$, $F_{1,3} = 0.036$, $P = 0.861$), or sham-treated RTT mice (RTT-sham; $n = 5$, $F_{1,4} = 0.018$, $P = 0.901$). DBS reduced the amplitude of the evoked population spikes from the baseline test in DBS-treated RTT mice (RTT-DBS; $n = 5$, $F_{1,4} = 6.73$, $P = 0.060$). * $P < 0.05$ (Tukey's post hoc). All data are presented as mean \pm s.e.m.



Extended Data Figure 7 | Unilateral fornical DBS induces neuronal activity and stimulates neurogenesis bilaterally in the dentate gyrus.
a, Representative images showing that expression of the *Fos* gene was increased following fornical DBS in wild-type and RTT mice compared to their sham controls, respectively (percentage of ipsilateral c-Fos-positive cells over the dentate granule cells: WT-sham, $0.26 \pm 0.04\%$; WT-DBS, $34.52 \pm 4.62\%$;

RTT-sham, $0.30 \pm 0.05\%$; RTT-DBS, $32.55 \pm 3.74\%$). **b**, Representative images showing that there were more BrdU⁺ (green), DCX⁺ (red), and merged (yellow) cells in the dentate gyrus in fornical DBS-treated wild-type and RTT mice than in their respective sham controls. Scale bar, 100 μ m. Con, contralateral; Ips, ipsilateral.



Extended Data Figure 8 | The cholinergic antagonist atropine did not alter fornical DBS-induced enhancement of fear memory. **a**, Placement of guide cannula and recording electrode into the dorsal hippocampus. **b**, Hippocampal infusion of 1.0 μ g atropine did not change the amplitudes of the evoked potentials of the Ffx recorded in the dentate gyrus in both RTT and wild-type mice. There was no difference of the population spike amplitudes before or after atropine infusion in both RTT mice ($n = 5$; one-way ANOVA, $F_{9,36} = 0.69$, $P = 0.715$) and wild-type controls ($n = 3$; $F_{9,18} = 0.99$, $P = 0.485$). **c**, Representative hippocampal EEG traces before and after vehicle (V) or atropine (A) infusion. Scale bars: 0.5 s, 0.2 mV. **d**, RTT mice ($n = 17$) showed less spontaneous hippocampal theta activity than wild-type animals ($n = 20$) (** $P < 0.01$, two-tailed t -test). **e**, Hippocampal infusion of atropine, but not vehicle, reduced hippocampal theta oscillation in both RTT and wild-type mice compared to their pre-infusion baselines (WT-V, $n = 9$; WT-A, $n = 11$; RTT-V, $n = 8$; RTT-A, $n = 9$; * $P < 0.05$, two-tailed paired t -test; n.s., not significant). **f**, Hippocampal microinfusion of atropine before fear conditioning training did not alter fear memory in fornical DBS treated RTT mice or wild-type controls. Mice in all four groups (WT-V, $n = 10$; WT-A, $n = 11$; RTT-V, $n = 12$; RTT-A, $n = 13$) experienced 2 weeks of fornical DBS that was

finished 3 weeks before fear conditioning training. Atropine or vehicle was bilaterally infused into the dorsal hippocampus before training. Memory retention was tested 24 h after training. Two-way ANOVA revealed a significant main effect of genotype ($F_{1,42} = 10.27$, $P = 0.003$), but there was no difference between atropine- and vehicle-treated mice (treatment, $F_{1,42} = 0.34$, $P = 0.562$; genotype \times treatment interaction, $F_{1,42} = 0.069$, $P = 0.794$). Atropine did not change cued fear memory, either: two-way ANOVA revealed no difference between genotypes ($F_{1,42} = 2.99$, $P = 0.091$) or between atropine- and vehicle-treated mice (treatment, $F_{1,42} = 0.046$, $P = 0.831$; genotype \times treatment interaction, $F_{1,42} = 0.154$, $P = 0.697$). * $P < 0.05$; n.s., not significant (Tukey's post hoc). **g**, Intra-hippocampal atropine infusion alone did not change the basal level of freezing in the contextual test environment in either wild-type or RTT mice. There was no difference between vehicle- ($n = 9$) or atropine-treated ($n = 6$) mice ($P > 0.05$, two-tailed t -test). **h**, Schematic representation of the dorsal hippocampus at seven rostral-caudal planes (according to ref. 31) for the microinfusion sites in DBS-treatment experiments. The numbers on the left represent the posterior coordinate from the bregma. All data are presented as mean \pm s.e.m.

Control of REM sleep by ventral medulla GABAergic neurons

Franz Weber¹, Shinjae Chung¹, Kevin T. Beier², Min Xu¹, Liqun Luo² & Yang Dan¹

Rapid eye movement (REM) sleep is a distinct brain state characterized by activated electroencephalogram and complete skeletal muscle paralysis, and is associated with vivid dreams^{1–3}. Transection studies by Jouvet first demonstrated that the brainstem is both necessary and sufficient for REM sleep generation², and the neural circuits in the pons have since been studied extensively^{4–8}. The medulla also contains neurons that are active during REM sleep^{9–13}, but whether they play a causal role in REM sleep generation remains unclear. Here we show that a GABAergic (γ -aminobutyric-acid-releasing) pathway originating from the ventral medulla powerfully promotes REM sleep in mice. Optogenetic activation of ventral medulla GABAergic neurons rapidly and reliably initiated REM sleep episodes and prolonged their durations, whereas inactivating these neurons had the opposite effects. Optrode recordings from channelrhodopsin-2-tagged ventral medulla GABAergic neurons showed that they were most active during REM sleep (REM_{max}), and during wakefulness they were preferentially active during eating and grooming. Furthermore, dual retrograde tracing showed that the rostral projections to the pons and midbrain and caudal projections to the spinal cord originate from separate ventral medulla neuron populations. Activating the rostral GABAergic projections was sufficient for both the induction and maintenance of REM sleep, which are probably mediated in part by inhibition of REM-suppressing GABAergic neurons in the ventrolateral periaqueductal grey. These results identify a key component of the pontomedullary network controlling REM sleep. The capability to induce REM sleep on command may offer a powerful tool for investigating its functions.

Previous studies showed that the ventral medulla (vM) contains GABAergic neurons expressing the immediate early gene *c-fos* after deprivation-induced REM sleep rebound^{11,12}, suggesting REM sleep-related activity. To test the causal relationship between vM GABAergic neuron activity and brain states, we injected Cre-inducible adeno-associated viruses (AAV) expressing channelrhodopsin 2 fused with enhanced yellow fluorescent protein (ChR2-eYFP) locally into the vM of GAD2-Cre mice (Fig. 1a). Laser stimulation (20 Hz, 120 s per trial) was applied every 15–25 min, and brain states—wake, REM, and non-REM (NREM) sleep—were classified on the basis of electroencephalogram (EEG) and electromyogram (EMG) recordings. REM sleep was observed at a high probability during laser stimulation (Fig. 1b and Supplementary Video 1). To quantify the effect, we aligned all trials from six mice by the time of laser stimulation (Fig. 1c). We found a rapid, ~12-fold increase in REM sleep within 30 s of laser onset ($P < 0.001$, bootstrap) and a complementary decrease of NREM sleep. The EEG power spectrum and EMG activity during the laser-induced REM state were indistinguishable from those during spontaneous REM sleep outside laser stimulation periods (Extended Data Fig. 1). In control mice expressing eYFP without ChR2, laser had no effect (Extended Data Fig. 2), and the laser-induced change in the probability of REM sleep was significantly different between the ChR2 and control mice ($P < 0.001$, bootstrap). Furthermore, ChR2-mediated activation

of vM glutamatergic neurons reliably induced wakefulness rather than REM sleep (Extended Data Fig. 3), indicating that the REM-promoting effect was specific to GABAergic neurons.

The complementary changes in the probabilities of REM and NREM (Fig. 1c) suggest that vM GABAergic neuron activation primarily triggered NREM to REM transitions. To test this possibility, we analysed the effect of laser on the transition probability between each pair of brain states. Laser stimulation markedly enhanced the NREM to REM transitions ($P = 0.02$, bootstrap; Extended Data Fig. 4a). As a result, for trials in which laser onset fell on NREM sleep, the probability

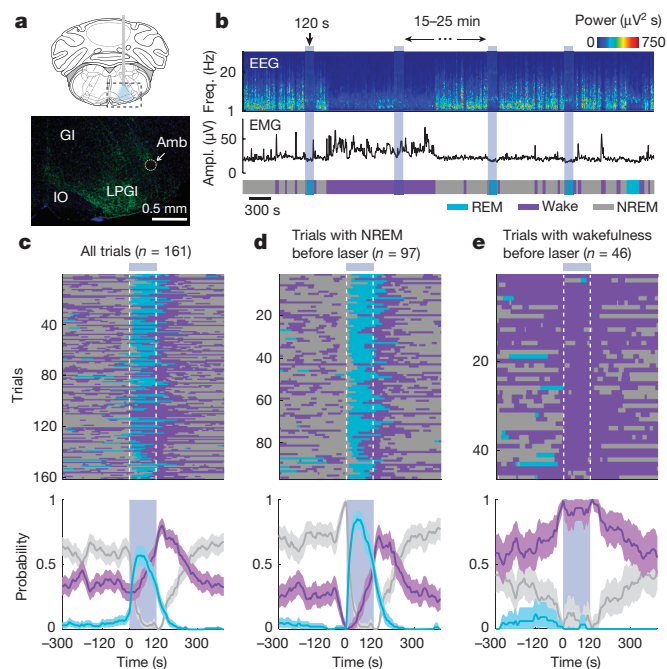


Figure 1 | Optogenetic activation of vM GABAergic neurons induces REM sleep. **a**, Top, schematic of optogenetic experiment (adapted with permission from The Mouse Brain in Stereotaxic Coordinates, 3rd edn, Franklin, K.B.J. and Paxinos, G., 88, Copyright Elsevier (Academic Press, 2007)³¹). Bottom, fluorescence image of vM (dashed box in schematic) in a GAD2-Cre mouse injected with AAV expressing ChR2-eYFP (green). Blue, 4',6-diamidino-2-phenylindole (DAPI). GI, gigantocellular reticular nucleus; LPGI, lateral paragigantocellular nucleus; IO, inferior olive; Amb, ambiguous nucleus. ChR2-eYFP was expressed within ~400 μ m from injection site and mainly localized within the LPGI. **b**, Example experiment. Shown are EEG power spectrogram (Freq., frequency), EMG amplitude (Ampl.), and brain states (colour coded). Blue shading, laser stimulation period (20 Hz, 120 s). **c**, Brain states in all trials from six mice (top) and probability of wake, NREM, or REM states (bottom) before, during, and after laser stimulation. Shading, 95% CI. Blue bar, laser stimulation period. **d**, Similar to c, for trials with laser onset falling on NREM sleep (probability of NREM is 1 immediately before laser onset). **e**, Trials with laser onset falling on wakefulness.

¹Division of Neurobiology, Department of Molecular and Cell Biology, Helen Wills Neuroscience Institute, Howard Hughes Medical Institute, University of California, Berkeley, California 94720, USA.

²Department of Biology, Howard Hughes Medical Institute, Stanford University, Stanford, California 94305, USA.

of REM sleep increased to 85% within 50 s (Fig. 1d), and NREM to REM transitions occurred in 94% of these trials. Notably, laser stimulation caused no significant change in NREM to wake ($P = 0.24$) and a reduction in the probability of REM to wake ($P < 0.001$) transitions (Extended Data Fig. 4b, c). Thus the gradual increase of wakefulness during laser stimulation (Fig. 1c, purple) was not due to a direct induction of wakefulness, but was a secondary consequence of laser-induced increase in REM sleep (in rodents, periods of REM sleep are often shorter than 120 s and typically followed by wakefulness). Furthermore, direct wake to REM transitions, characteristic of narcolepsy, were not observed (Extended Data Fig. 4e), indicating that optogenetic activation of vM GABAergic neurons triggered normal but not pathological transitions into REM sleep. This also explains why in trials in which the laser was turned on during wakefulness, no increase in REM sleep was observed (Fig. 1e and trial 3 in Supplementary Video 1).

To test whether these neurons also contribute to the maintenance of REM sleep, we applied a closed-loop protocol¹⁴, in which laser stimulation was initiated after spontaneous onset of REM sleep and maintained throughout the REM episode (Fig. 2a). The stimulation was applied randomly in ~50% of the episodes, and REM duration was compared between the episodes with and without stimulation (Fig. 2b). Laser stimulation increased the mean duration of REM episodes by 111% in GAD2-ChR2 mice (laser, 148.0 ± 19.2 s (mean \pm s.d.); no laser, 70.2 ± 10.1 s; $P = 0.00037$, paired t -test, $n = 6$ mice) but not in eYFP control mice (laser, 77.0 ± 14.9 s; no laser, 83.8 ± 13.6 s; $P = 0.62$, $n = 4$; Fig. 2c). Conversely, archaerhodopsin (Arch)- or halorhodopsin (Halo)-mediated silencing of vM GABAergic neurons (green laser) caused a 41% reduction of REM duration (laser, 54.1 ± 9.8 s; no laser, 91.2 ± 11.7 s; $P = 0.001$, $n = 4$), an effect not observed in eYFP control mice (laser, 90.0 ± 12.9 s; no laser, 79.5 ± 9.6 s; $P = 0.075$, $n = 4$; Fig. 2d). Thus, vM GABAergic neuron activity also contributes strongly to REM sleep maintenance. In addition, we pharmacogenetically silenced vM GABAergic neurons by expressing hM4Di¹⁵. Compared with vehicle injection, inhibition of these neurons by clozapine-*N*-oxide (CNO) injection strongly reduced REM sleep in a dose-dependent manner (Extended Data Fig. 5), further indicating the importance of vM GABAergic neurons in REM sleep generation.

Although vM GABAergic neurons have been shown to express *c-fos* after deprivation-induced REM sleep rebound^{11,12}, the slow time course of *c-fos* expression and its modulation by non-activity-related factors¹⁶ limit its precision in identifying REM-active neurons. To

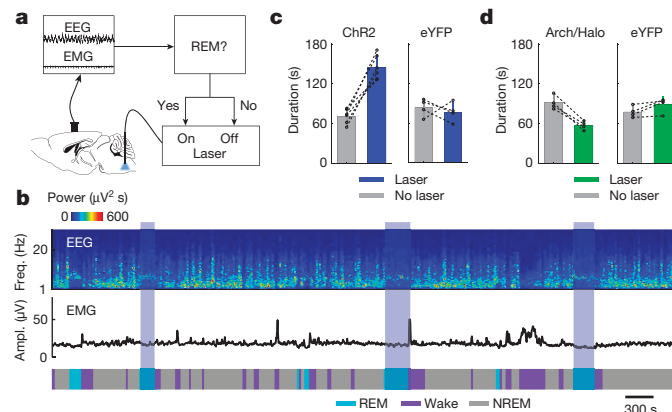


Figure 2 | Activation of vM GABAergic neurons prolongs REM sleep duration. **a**, Schematic of closed-loop stimulation. Laser was turned on after spontaneous REM onset and turned off at termination of the REM episode. **b**, Example recording containing REM episodes with and without laser stimulation (blue shading). **c**, Mean REM sleep duration with and without laser stimulation, in ChR2 (left, $n = 6$) and eYFP control (right, $n = 4$) mice. Each pair of dots, data from one mouse. Error bar, s.d. **d**, Similar to **c**, but with green laser stimulation in Arch/Halo (left, $n = 4$) and eYFP control (right, $n = 4$) mice.

understand how vM GABAergic neurons regulate REM sleep under natural conditions, we recorded their spiking activity across brain states. Since the vM contains multiple spatially intermingled cell types¹⁷, we tagged the GABAergic neurons with ChR2 by crossing GAD2-Cre with ChR2 reporter (Ai 32) mice. Recordings were made in freely moving mice by using optrodes, consisting of an optical fibre surrounded by multiple microelectrodes. High-frequency laser pulses (15 or 30 Hz) were applied intermittently, and single units exhibiting reliable laser-evoked responses with short latencies and low jitter were identified as GABAergic neurons (Fig. 3a–c and Extended Data Fig. 6a). The identified neurons consistently fired at high rates during REM sleep, low rates during NREM sleep, and variable rates during wakefulness⁹ (Fig. 3d and Supplementary Video 2). When the firing rates were averaged within each state, all 20 identified GABAergic neurons were most active during REM sleep ($P < 0.005$, Wilcoxon rank-sum test; Fig. 3e–g, referred to as ‘REM_{max}’ neurons). The mean firing rate of these neurons during REM sleep (34.6 spikes per second, 95% confidence interval (CI) 23.6–48.0 spikes per second) was in fact higher than the laser stimulation frequency applied in our experiments (Figs 1 and 2), suggesting that their activity is sufficient for the observed REM promoting effects. Their firing rates increased gradually over ~30 s before the NREM to REM transition and decreased abruptly at the end of REM sleep (Extended Data Fig. 7). Such a temporal profile is well suited for both the induction and maintenance of REM sleep.

Since during wakefulness the firing rates of vM GABAergic neurons were highly variable, we analysed their relationship with different

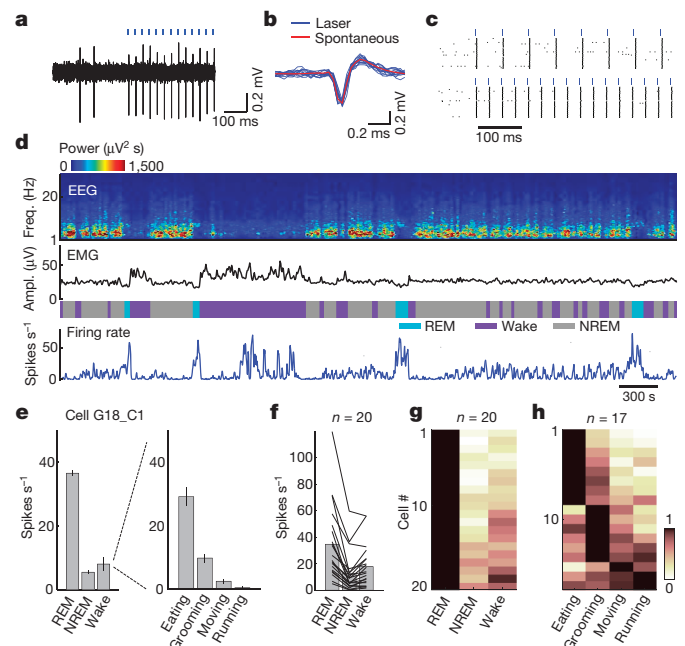


Figure 3 | Firing rates of identified vM GABAergic neurons across brain states. **a**, Example recording of spontaneous and laser-evoked spikes from a vM neuron. Blue ticks, laser pulses (30 Hz). **b**, Comparison between laser-evoked (blue) and averaged spontaneous (red) spike waveforms from this unit. **c**, Spike raster showing multiple trials of laser stimulation at 15 and 30 Hz. **d**, Firing rates of an example vM GABAergic neuron. **e**, Left, mean firing rates of the neuron in **d** during different brain states. Error bar, s.e.m. Right, firing rates of the neuron during different wakeful behaviours. **f**, Firing rates of 20 identified vM GABAergic neurons during different brain states. Each line shows firing rates of one unit; grey bar, average across units. **g**, Relative firing rates across brain states. The rates of each neuron were normalized by its maximum. All 20 neurons showed maximum firing rates during REM sleep (the difference was significant at $P < 0.005$ for all 20 neurons compared with NREM and wakefulness; Wilcoxon rank-sum test, post-hoc Bonferroni correction). **h**, Relative firing rates of 17 vM GABAergic neurons during different wakeful behaviours.

wakeful behaviours (classified on the basis of video recordings (Supplementary Video 2)). Of the 17 identified GABAergic neurons with recording periods encompassing multiple behaviours, 82% were most active during eating ($n = 8$) or grooming ($n = 6$) and much less active during moving or running (Fig. 3h and Extended Data Fig. 6c). Such a firing pattern is distinct from that of noradrenergic neurons in the locus coeruleus (a known postsynaptic target of vM GABAergic neurons¹⁸), which are strongly suppressed during eating and grooming¹⁹. Interestingly, activation of vM GABAergic neurons during wakefulness increased eating (Extended Data Fig. 8), indicating a causal relationship between their firing and the eating behaviour. Among the 24 unidentified neurons recorded in the vM (likely to include non-GABAergic neurons), 71% were also REM_{max}^{9,10,13} (Extended Data Fig. 9a, b). Compared with the identified GABAergic neurons, however, the unidentified group showed a delayed firing rate increase before REM sleep onset (~ 10 s) and a more gradual decrease at the end of REM sleep (Extended Data Fig. 9c, d). In addition, many of them were most active during running or moving (Extended Data Fig. 9b, right).

Previous studies showed that vM neurons project both rostrally to the pons, midbrain, and other brain regions and caudally to the spinal cord^{13,20}. We wondered which pathway mediates the REM-promoting effect. The pons and midbrain are known to be crucial for REM sleep generation^{2,4,6–8,12}. Injection of Cre-inducible AAV expressing ChR2-eYFP into the vM of GAD2-Cre mice revealed extensive axon projections to these regions (Fig. 4a). In particular, the ventrolateral periaqueductal grey (vIPAG) is thought to suppress REM sleep by providing GABAergic inhibition to REM-promoting neurons^{7,12}. To test the effect of vIPAG GABAergic neuron activity on brain states, we injected Cre-inducible AAV expressing ChR2-eYFP into the vIPAG of GAD2-Cre mice. Laser stimulation almost completely suppressed onset of REM sleep ($P < 0.001$, bootstrap, $n = 6$) and shortened the duration of REM episodes (laser, 37.1 ± 10.2 s; no laser, 82.8 ± 9.9 s; $P = 0.003$, paired t -test, $n = 5$) while strongly promoting NREM sleep ($P < 0.001$, bootstrap; Fig. 4c). Rabies-virus-mediated monosynaptic retrograde tracing confirmed that vIPAG GABAergic neurons receive direct inhibitory innervation from vM neurons (Fig. 4d, e). Furthermore, optogenetic activation of vM GABAergic neuron axons within the vIPAG caused large increases both in the initiation ($P < 0.001$, bootstrap, $n = 5$) and in the duration (laser, 130.7 ± 32.9 s; no laser, 59.6 ± 9.6 s; $P = 0.02$, paired t -test, $n = 5$) of REM sleep episodes (Fig. 4f), with magnitudes comparable to those found with stimulation of vM cell bodies (Figs 1 and 2). Although in principle stimulation of axon fibres in the vIPAG could antidromically activate neuronal cell bodies in the vM and axon collaterals to other regions, simultaneous injections of retrograde tracers to both the pons and spinal cord showed little overlap between the labelled neurons (1.3%, Fig. 4b), indicating that the rostral and caudal projections arise largely from separate vM neuron populations. Thus, the rostral projections from vM GABAergic neurons exert a strong REM sleep-promoting effect, probably mediated at least in part by direct inhibition of the REM-suppressing vIPAG GABAergic neurons.

Although the medulla is known to be involved in REM sleep generation, previous studies have focused primarily on its role in inducing muscle atonia, through projections to the spinal cord^{21,22}. Our findings indicate that vM GABAergic neurons constitute a critical component of the core network, generating not only muscle atonia but also the cortical activation associated with REM sleep^{17,23,24}. The effect is largely, if not completely, attributable to the rostral projections (Fig. 4), consistent with the previous finding that inactivation of the pons or transection at the pontomedullary junction blocks the muscle atonia evoked by stimulation of the medulla^{25,26}. The projection to the locus coeruleus is likely to provide GABAergic inhibition of noradrenergic neurons^{13,18}, whose activity enhances wakefulness²⁷ and excitability of spinal motor neurons²⁸. Inhibition of the vIPAG GABAergic neurons (Fig. 4c) should cause disinhibition of the REM-promoting neurons in the pons^{5–8,12}, which may in turn trigger muscle atonia

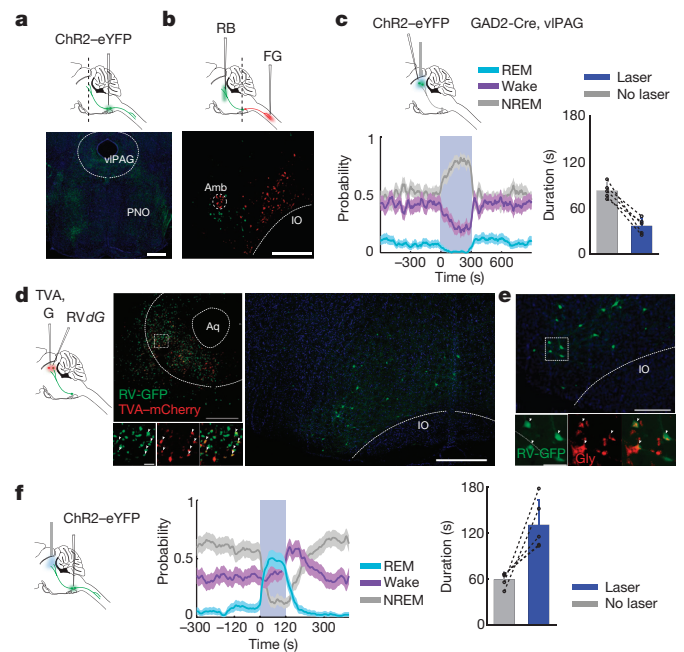


Figure 4 | Inhibition of vIPAG GABAergic neurons by vM projections promotes REM sleep. **a**, Top, schematic showing injection of AAV expressing ChR2-eYFP into the vM of a GAD2-Cre mouse. Bottom, fluorescence image of vM axons in the pons and midbrain (position of coronal section indicated by dashed line in schematic). Scale bar, 500 μ m; green, eYFP; blue, DAPI. **b**, Top, schematic showing simultaneous injections of RetroBeads (RB) in the pons and FluoroGold (FG) in the spinal cord. Bottom, fluorescence image of vM, showing neurons labelled by RB and FG. Among the 1235 FG- and 881 RB-labelled neurons, only 27 were double labelled. Scale bar, 500 μ m. **c**, Left, probability of wake, NREM, or REM states before, during, and after laser stimulation of vIPAG GABAergic neurons (20 Hz, 300 s; $n = 6$ mice). Shading, 95% CI; Blue bar, laser stimulation period (20 Hz, 300 s). Right, mean REM durations with and without vIPAG stimulation ($n = 5$ mice). Each pair of dots, data from one mouse. Error bar, s.d. **d**, Left, schematic showing rabies-mediated trans-synaptic tracing. TVA, EnvA receptor; G, rabies glycoprotein; RVdG, G-deleted rabies virus. Middle, fluorescence image of vIPAG in a GAD2-Cre mouse. Scale bar, 500 μ m. Bottom middle, enlarged view of region in white box showing starter cells (yellow, expressing both GFP and mCherry, arrowheads; scale bar, 20 μ m). Right, rabies-labelled presynaptic neurons in vM (same brain as in middle panel). Scale bar, 500 μ m. **e**, Rabies-labelled presynaptic neurons in vM are GABAergic and glycinergic. Lower panel, enlarged view of region in white box, containing GFP-labelled neurons expressing glycine (Gly, arrowheads; scale bar, 50 μ m). GABA and glycine coexist in a high percentage of vM neurons (Extended Data Fig. 10), suggesting that the glycinergic rabies-labelled neurons in the vM are also GABAergic. In total, 82% (185/226) rabies-labelled cells were glycine positive ($n = 3$ mice). **f**, Left, probability of wake, NREM, or REM states before, during, and after laser stimulation of vM axons in vIPAG (20 Hz, 120 s; $n = 5$ mice). Shading, 95% CI. Right, mean REM durations with and without vM axon stimulation ($n = 5$ mice). Error bar, s.d.

through their projections to the medulla and spinal cord, and cortical activation through their projections to midbrain and forebrain regions. The glutamatergic REM-active neurons in the dorsal pons⁶ may also innervate the vM GABAergic neurons¹³, thus shaping their firing rates across different brain states (Fig. 3).

In addition to the vM, the dorsal medulla also contains REM_{max} GABAergic neurons projecting to the locus coeruleus²⁹, which may also contribute to the generation of REM sleep. The induction of NREM to REM, but not wake to REM, transitions by vM neurons points to a robust mechanism of wake maintenance—probably involving orexin/hypocretin neurons³⁰—that cannot be overridden by activating the vM GABAergic neurons. How this mechanism interacts with the circuit that generates REM sleep remains to be elucidated. Furthermore, a long-standing mystery is what functions are served by

REM sleep and its associated dreaming. The ability to control REM sleep at a high temporal precision, as demonstrated in this study, provides a powerful tool for studying its functions.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 14 May; accepted 23 June 2015.

Published online 7 October 2015.

1. Aserinsky, E. & Kleitman, N. Regularly occurring periods of eye motility, and concomitant phenomena, during sleep. *Science* **118**, 273–274 (1953).
2. Jouvet, M. Recherches sur les structures nerveuses et les mécanismes responsables des différentes phases du sommeil physiologique. *Arch. Ital. Biol.* **100**, 125–206 (1962).
3. Dement, W. The occurrence of low voltage, fast, electroencephalogram patterns during behavioral sleep in the cat. *Electroencephalogr. Clin. Neurophysiol.* **10**, 291–296 (1958).
4. Hobson, J. A., McCarley, R. W. & Wyzinski, P. W. Sleep cycle oscillation: reciprocal discharge by two brainstem neuronal groups. *Science* **189**, 55–58 (1975).
5. Boissard, R., Fort, P., Gervasoni, D., Barbagli, B. & Luppi, P. H. Localization of the GABAergic and non-GABAergic neurons projecting to the sublaterodorsal nucleus and potentially gating paradoxical sleep onset. *Eur. J. Neurosci.* **18**, 1627–1639 (2003).
6. Clément, O., Sapin, E., Bédard, A., Fort, P. & Luppi, P. H. Evidence that neurons of the sublaterodorsal tegmental nucleus triggering paradoxical (REM) sleep are glutamatergic. *Sleep* **34**, 419–423 (2011).
7. Lu, J., Sherman, D., Devor, M. & Saper, C. B. A putative flip-flop switch for control of REM sleep. *Nature* **441**, 589–594 (2006).
8. Van Dort, C. J. *et al.* Optogenetic activation of cholinergic neurons in the PPT or LDT induces REM sleep. *Proc. Natl Acad. Sci. USA* **112**, 584–589 (2015).
9. Siegel, J. M., Wheeler, R. L. & McGinty, D. J. Activity of medullary reticular formation neurons in the unrestrained cat during waking and sleep. *Brain Res.* **179**, 49–60 (1979).
10. Sakai, K., Kanamori, N. & Jouvet, M. Activités unitaires spécifiques du sommeil paradoxal dans la formation réticulée bulbaire chez le chat non-restreint. *C.R. Seances Acad. Sci. D* **289**, 557–561 (1979).
11. Maloney, K. J., Mainville, L. & Jones, B. E. c-Fos expression in GABAergic, serotonergic, and other neurons of the pontomedullary reticular formation and raphe after paradoxical sleep deprivation and recovery. *J. Neurosci.* **20**, 4669–4679 (2000).
12. Sapin, E. *et al.* Localization of the brainstem GABAergic neurons controlling paradoxical (REM) sleep. *PLoS One* **4**, e2272 (2009).
13. Sirieix, C., Gervasoni, D., Luppi, P. H. & Léger, L. Role of the lateral parabrachial nucleus in the network of paradoxical (REM) sleep: an electrophysiological and anatomical study in the rat. *PLoS One* **7**, e28724 (2012).
14. Jégo, S. *et al.* Optogenetic identification of a rapid eye movement sleep modulatory circuit in the hypothalamus. *Nature Neurosci.* **16**, 1637–1643 (2013).
15. Armbruster, B. N., Li, X., Pausch, M. H., Herlitze, S. & Roth, B. L. Evolving the lock to fit the key to create a family of G protein-coupled receptors potentially activated by an inert ligand. *Proc. Natl Acad. Sci. USA* **104**, 5163–5168 (2007).
16. Kovács, K. J. c-Fos as a transcription factor: a stressful (re)view from a functional map. *Neurochem. Int.* **33**, 287–297 (1998).
17. Holmes, C. J. & Jones, B. E. Importance of cholinergic, GABAergic, serotonergic and other neurons in the medial medullary reticular formation for sleep-wake states studied by cytotoxic lesions in the cat. *Neuroscience* **62**, 1179–1200 (1994).
18. Aston-Jones, G., Ennis, M., Pieribone, V. A., Nickell, W. T. & Shipley, M. T. The brain nucleus locus coeruleus: restricted afferent control of a broad efferent network. *Science* **234**, 734–737 (1986).
19. Aston-Jones, G. & Bloom, F. E. Activity of norepinephrine-containing locus coeruleus neurons in behaving rats anticipates fluctuations in the sleep-waking cycle. *J. Neurosci.* **1**, 876–886 (1981).
20. Loewy, A. D., Wallach, J. H. & McKellar, S. Efferent connections of the ventral medulla oblongata in the rat. *Brain Res. Rev.* **3**, 63–80 (1981).
21. Magoun, H. W. & Rhines, R. An inhibitory mechanism in the bulbar reticular formation. *J. Neurophysiol.* **9**, 165–171 (1946).
22. Schenkel, E. & Siegel, J. M. REM sleep without atonia after lesions of the medial medulla. *Neurosci. Lett.* **98**, 159–165 (1989).
23. Vanni-Mercier, G., Sakai, K., Lin, J. S. & Jouvet, M. Carbachol microinjections in the mediodorsal pontine tegmentum are unable to induce paradoxical sleep after caudal pontine and prebulbar transections in the cat. *Neurosci. Lett.* **130**, 41–45 (1991).
24. Webster, H. H., Friedman, L. & Jones, B. E. Modification of paradoxical sleep following transections of the reticular formation at the pontomedullary junction. *Sleep* **9**, 1–23 (1986).
25. Kohyama, J., Lai, Y. Y. & Siegel, J. M. Inactivation of the pons blocks medullary-induced muscle tone suppression in the decerebrate cat. *Sleep* **21**, 695–699 (1998).
26. Siegel, J. M., Nienhuis, R. & Tomaszewski, K. S. Rostral brainstem contributes to medullary inhibition of muscle tone. *Brain Res.* **268**, 344–348 (1983).
27. Carter, M. E. *et al.* Tuning arousal with optogenetic modulation of locus coeruleus neurons. *Nature Neurosci.* **13**, 1526–1533 (2010).
28. White, S. R., Fung, S. J. & Barnes, C. D. Norepinephrine effects on spinal motoneurons. *Prog. Brain Res.* **88**, 343–350 (1991).
29. Kaur, S., Saxena, R. N. & Mallick, B. N. GABAergic neurons in prepositus hypoglossi regulate REM sleep by its action on locus coeruleus in freely moving rats. *Synapse* **42**, 141–150 (2001).
30. Taheri, S., Zeitzer, J. M. & Mignot, E. The role of hypocretins (orexins) in sleep regulation and narcolepsy. *Annu. Rev. Neurosci.* **25**, 283–313 (2002).
31. Franklin, K. B. J. & Paxinos, G. *The Mouse Brain in Stereotaxic Coordinates* 3rd edn, 88 (Academic Press, 2007).

Supplementary Information is available in the online version of the paper.

Acknowledgements We thank A. Popescu for the help with *in vivo* physiology, M. Bikov and S. Chung for technical assistance, the University of North Carolina Virus Core for supplying AAV, and T. Kilduff and J. Cox for discussions. This work was supported by EMBO and Human Frontier Science Program postdoctoral fellowships (to F.W.).

Author Contributions F.W. and Y.D. conceived and designed the experiments. F.W. performed all optogenetic stimulation experiments and optrode recordings. S.C. performed a subset of pharmacogenetic experiments and fluorescence microscopy. K.T.B. and L.L. provided viral reagents for rabies-mediated trans-synaptic experiments. M.X. designed the optrodes used in this study. F.W. and Y.D. wrote the manuscript, and all authors participated in the revision of the manuscript.

Author Information All primary histological, electrophysiological, and behavioural data have been archived in the Department of Molecular and Cell Biology, University of California, Berkeley. Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to Y.D. (ydan@berkeley.edu).

METHODS

Animals. All experimental procedures were approved by the Animal Care and Use Committee at the University of California, Berkeley. Optogenetic manipulation experiments were performed in male or female GAD2-Cre (Jackson Laboratory stock 010702) and VGLUT2-Cre (016963) mice. Pharmacogenetic, rabies-mediated trans-synaptic tracing and anterograde tracing experiments were performed in male or female GAD2-Cre mice. For optrode recording experiments, male GAD2-Cre mice were crossed with female loxP-flanked-ChR2-eYFP mice (012569). Retrograde tracing experiments with RetroBeads (Lumafuor) and FluoroGold (Fluorochrome) were performed in male or female wild-type (C57) mice.

Surgery. Adult (6- to 12-week-old) mice were anaesthetized with isoflurane (3% induction, 1.5% maintenance) and placed on a stereotaxic frame. Body temperature was kept stable throughout the procedure by using a heating pad. After asepsis, the skin was incised to expose the skull, and the overlying connective tissue was removed. For optogenetic experiments, a craniotomy (~1 mm diameter) was made above the right cerebellum (6.8 mm posterior to the bregma, 0.6–0.9 mm lateral). Virus (AAV2-EF1 α -FLEX-ChR2-eYFP, AAV2-EF1 α -FLEX-eYFP, AAV2-EF1 α -FLEX-eNpHR3.0-eYFP, AAV2-CAG-FLEX-ArchT-tdTomato, produced by University of North Carolina Vector Core, or AAV-DJ-EF1 α -FLEX-ChR2-eYFP) was loaded into a sharp micropipette mounted on a Nanoject II and injected slowly at a depth of 4.8 mm from the brain surface (600–800 nl). An optical fibre (200 μ m diameter) was inserted with the tip ~400 μ m above the virus injection site. To stimulate axon projections of vM GABAergic neurons in the vPAG, the optical fibre was implanted on top of the vPAG (4.8 mm posterior to the bregma, 0.6 mm lateral, 2.3 mm depth). To activate vPAG GABAergic neurons, AAV expressing ChR2-eYFP was injected into the vPAG. For pharmacogenetic experiments, AAV8-hSyn-FLEX-hM4D(Gi)-mCherry was injected bilaterally (300–800 nl on each side). For EEG and EMG recordings, a reference screw was inserted into the skull on top of the left cerebellum. EEG recordings were made from two screws on top of the left and right cortex (1 mm posterior to the bregma). Two EMG electrodes were inserted into the left and right neck muscles. The optical fibre, screws, and EEG/EMG electrodes were secured to the skull using dental cement. Without optogenetic manipulation, the sleep architecture quantified by the percentage, average duration and frequency of each brain state of these implanted animals was well within the range reported for mice without these implants, indicating that virus injection and optical fibre implantation did not cause large chronic changes in sleep architecture. For optogenetic experiments, data from animals where the tip of the optical fibre was not within the aimed location were excluded. For pharmacogenetic experiments, data from animals where virus expression was not restricted to the vM were excluded.

For optrode recording experiments, the optrode assembly was inserted at a depth of 4.8 mm. A screw was attached to the frontal skull for EEG recording, and an EMG electrode was inserted into the right neck muscle. The optrode assembly, screws, and EEG/EMG electrodes were secured to the skull using dental cement.

For rabies-mediated retrograde tracing experiments, we first injected 100–150 nl of a combination of AAV5-CAG-FLEX-TC^B (2.6×10^{12} vector genomes per millilitre) and AAV8-CAG-FLEX-RG (1.3×10^{12} vector genomes per millilitre)³² (addgene 48332 and 48333, respectively) into the vPAG (4.8 mm posterior to the bregma, 0.6 mm lateral, 2.5 mm depth) of GAD2-Cre mice. Three weeks later, RVdG (titre 5×10^8 colony forming units per millilitre) pseudotyped with EnvA was injected into the vPAG. AAV5-CAG-FLEX-TC^B and AAV-CAG-FLEX-RG were purchased from the University of North Carolina viral core, and RVdG was amplified in house from a stock derived from the Salk Institute viral core. Data from mice where starter cells were found outside the vPAG were excluded. For controls, we injected the same viruses in the same temporal sequence into wild-type mice ($n = 2$). No GFP-labelled cells were detected in the vM.

For dual retrograde tracing experiments, we injected RetroBeads bilaterally into the pons (5.0–5.2 mm posterior to the bregma, 0.9 mm lateral) at several depths (3.0, 3.5, and 4.0 mm, 150–200 nl at each depth). The spinal cord was exposed by a unilateral laminectomy of a cervical vertebra. After a small incision in the dura, FluoroGold was injected (0.7 mm lateral, 1 mm depth, 400–700 nl). Histology and *in vivo* experiments were performed more than 10 and 21 days after injection.

Polysomnographic recordings. Animals were housed on a 12-h dark/12-h light cycle (light on between 7:00 and 19:00). Behavioural experiments were performed between 13:30 and 18:30. EEG and EMG electrodes were connected to flexible recording cables via a mini-connector, and recordings were made in the animal's home cage placed in a sound-attenuation box. Recordings started after more than 1 h of habituation. The signals were recorded with a TDT RZ5 amplifier, filtered (1–750 Hz) and digitized at 1,500 Hz. Brain states were classified into NREM sleep, REM sleep, and wakefulness using custom-written MATLAB software. The classification was performed without any information about the identity of the animal or laser stimulation timing. First, we calculated the power spectrum of the EEG

and EMG using a 5 s sliding window, sequentially shifted by 2.5 s increments. Next, we summed the EEG power in the ranges from 1 to 4 Hz and from 6 to 12 Hz, yielding a time-dependent delta and theta power, respectively. For further analysis, we divided the theta by the delta power (theta/delta ratio). We also computed the total EMG power from 20 to 300 Hz. For each time point, we determined the brain state using a threshold algorithm. A state was classified as NREM if the delta power was lower than its mean (averaged over the whole recording session) and if the EMG power was lower than its mean plus one standard deviation. A state was classified as REM if (1) the delta power was lower than the average, (2) the theta/delta ratio deviated more than one standard deviation from its mean, and (3) the EMG power was lower than its mean plus one standard deviation. All remaining states were classified as wake. The wake state thus encompassed states with high EMG power (active awake), or low delta power without elevated EMG activity or theta/delta ratio (quiet awake). Finally, we manually verified the automatic classification to ensure that all states were correctly assigned.

Behavioural monitoring. To classify different wakeful behaviours, we made video recordings (sampling rate, 5 Hz) using a camera placed on top of the cage. Wakeful behaviours were divided into four categories: eating, grooming, moving, and running. The classification was performed manually, on the basis of the video, EEG, and EMG recordings, using a custom-written graphical user interface (programmed in MATLAB). The experimenter classifying the wakeful behaviours was blind to the timing of laser stimulation or the identity of the animal.

Optogenetic manipulation. To test the role of vM GABAergic neurons in REM sleep induction, we applied blue laser stimulation (473 nm, 5 mW at fibre tip) at 20 Hz (10 ms per pulse). Each trial lasted for 120 s (stimulation of vM GABAergic neurons and axon projections) or 300 s (stimulation of vPAG GABAergic neurons), and the inter-trial interval was chosen randomly from a uniform distribution between 15 and 25 min. For each animal we recorded at least two 5 h sessions.

To test the role of vM or vPAG GABAergic neurons in REM sleep maintenance, we applied a closed-loop stimulation protocol. The animal's brain state was classified by real-time analysis of EEG/EMG recordings. As soon as REM sleep was detected, the laser was turned on with 50% probability, and turned off only when the REM episode ended. This allowed comparison of the durations of REM episodes with and without laser stimulation within the same recording session. REM periods were assigned to either the experimental or control group according to a pseudo-random Bernoulli sequence generated before the start of the experiment. For each animal at least two 5 h sessions were recorded. In ChR2-mediated activation experiments, we used a blue laser (473 nm, 5 mW, 20 Hz). In Arch/Halo-mediated inhibition experiments, we used a green laser (532 nm, 20 mW, step pulse).

Pharmacogenetic manipulation. To inhibit vM GABAergic neurons, we injected CNO dissolved in 0.1 ml vehicle solution (PBS with 0.5% dimethylsulfoxide (DMSO)) into GAD2-Cre mice expressing hM4Di in the vM, 20 min before the recording session. CNO was administered intraperitoneally at two doses (1 and 5 mg kg⁻¹ body weight) on different days. For the control experiment, we injected the vehicle solution without CNO. For experimental and control recordings, animals were subjected to the same behavioural manipulations.

Optrode recording. To record the activity of vM GABAergic neurons, we used a custom-built optrode³³, consisting of an optical fibre (0.1 or 0.2 mm diameter) surrounded by 12 microwire electrodes (Stablohm 675) twisted into stereotrodes or tetrodes. The electrode tips were electroplated in a chloride-platinum solution to an impedance of ~600 k Ω . The optical fibre and electrodes were inserted into a screw-driven microdrive. The optrode was slowly lowered in 25–50 μ m steps to search for light-responsive neurons. The signals were recorded using a TDT RZ5 amplifier, filtered (0.3–8 kHz) and digitized at 25 kHz. Recordings were performed over a period of 1–2 months from each mouse. At the end of the experiment electrolytic lesions were made by passing a current (100 μ A, 10 s) through one or two electrodes to identify the end of the recording tract.

Spike sorting. Spikes were sorted offline on the basis of the waveform energy and the first three principal components of a spike waveform on each stereotrode or tetrode channel. Single units were identified either manually using the software Klusters (<http://neurosuite.sourceforge.net>) or automatically using the software KlustaKwik (<http://klustakwik.sourceforge.net>). The quality of each unit was assessed by the presence of a refractory period and quantified using isolation distance and L_{ratio} ³⁴. Units with an isolation distance <19 or $L_{\text{ratio}} > 0.1$ were discarded. For units recorded with stereotrodes, the median values of isolation distance and L_{ratio} were 34.4 and 0.001. For units recorded with tetrodes, the corresponding median values were 38.3 and 0.023.

Identification of GABAergic units. To identify ChR2-expressing vM GABAergic neurons, high-frequency laser pulse trains (15 and 30 Hz with duration of 1 and 0.5 s, respectively) were delivered every 1 or 2 min. The laser power and pulse

duration were optimized to identify light-responsive neurons without changing the brain state. A unit was identified as GABAergic if spikes were evoked by laser pulses (both 15 and 30 Hz trains) at high reliability (>0.6), short first-spike latency (<6 ms), and small jitter (<2 ms), and if the waveforms of the laser-evoked and spontaneous spikes were highly similar (correlation coefficient >0.9 ; Extended Data Fig. 6a). To compute the mean firing rate of each neuron in each brain state, spikes during the laser pulse trains were excluded. Twenty out of twenty-one identified GABAergic units had their recording periods encompassing all three brain states. For 17 units the recording periods encompassed all four wakeful behaviours.

In total we recorded from 21 identified GABAergic units and 24 unidentified units. However, the ratio (21/24) is probably affected by sampling bias and thus may not reflect the actual percentage of ChR2-labelled GABAergic neurons in the vM. In general, the success rate for finding a neuron that is reliably activated by laser is very low. To maximize the rate of data collection from identified neurons, we spent the 1–2 h on recording (instead of moving the electrode to a new location) only if we found at least one unit that appeared to be activated by the laser pulses. As a result, most of the unidentified units were recorded simultaneously with other laser-driven, putative GABAergic units.

Histology and immunohistochemistry. Mice were deeply anaesthetized and transcardially perfused with 0.1 M PBS followed by 4% paraformaldehyde (w/v) in PBS. For fixation, brains were kept overnight in 4% paraformaldehyde. For cryoprotection, brains were placed in 30% sucrose (w/v) in PBS solution for at least two nights. After embedding and freezing, brains were sectioned into 20, 30, or 60 μ m coronal slices using a cryostat. For immunohistochemistry, non-specific binding sites were blocked by incubating the sections in 2% donkey or goat serum in PBST (0.3% Triton X-100 in PBS). Brain sections were stained with an anti-choline acetyltransferase antibody (AB144, Millipore; 1:250) for cholinergic neurons, and anti-glycine (AB5020, Millipore; 1:100) for glycinergic neurons. To amplify the fluorescence of axon fibres expressing ChR2-eYFP, we used an antibody against GFP (A11122 or A11120, Life Technologies, 1:1,000). Brain sections were incubated with the primary antibody in blocking solution for two nights. Species-specific secondary antibodies conjugated with red or green Alexa fluorophores (donkey anti-mouse, A21202, 1:1,000; goat anti-rabbit, A11008, 1:1,000; donkey anti-goat, A11058, 1:250) in PBS were applied for 2 h. Finally, slides were washed for 2 h in PBS and mounted with Fluoromount. Fluorescence images were taken using a confocal microscope (LSM 710 AxioObserver Inverted 34-Channel Confocal) and Nanozoomer.

Statistics. For optogenetic experiments, GAD2-Cre mice were randomly assigned to control (injected with AAV expressing eYFP) and experimental groups (injected with AAV expressing ChR2-eYFP, eNpHR3.0-eYFP, or ArchT-

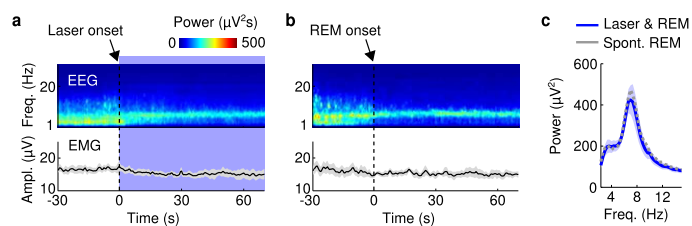
tdTomato). Experimental and control animals were subjected to exactly the same surgical and behavioural manipulations. No randomization was used for rabies-mediated, anterograde or dual retrograde tracing experiments, pharmacogenetic recordings and optrode recordings. Unless stated otherwise, investigators were not blinded to animal identity and outcome assessment.

All statistical tests (paired *t*-test, Wilcoxon rank-sum test, Wilcoxon signed-rank test, bootstrap) were two-sided. For both paired and unpaired tests, we ensured that the variances of the data were similar between the compared groups. For *t*-tests, we verified that the data were normally distributed using Lilliefors test for normality.

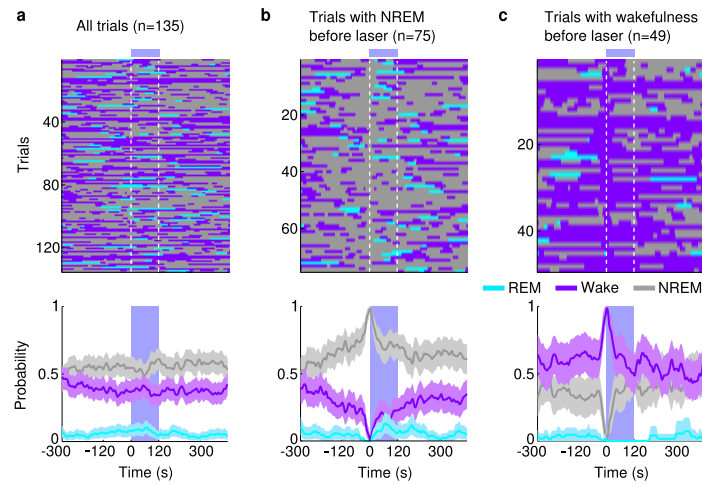
The 95% CIs for brain state probabilities were calculated using a bootstrap procedure. For an experimental group of *n* mice, with mouse *i* comprising *m_i* trials, we calculated the CI as follows: we repeatedly resampled the data by randomly drawing for each mouse *m_i* trials (random sampling with replacement). For each of the 10,000 iterations, we recalculated the mean probabilities for each brain state across the *n* mice. The lower and upper CIs were then extracted from the distribution of the resampled mean values. To test whether a given brain state was significantly modulated by laser stimulation, we calculated for each bootstrap iteration the difference between the mean probabilities during laser stimulation and the preceding period of identical duration. From the resulting distribution of difference values, we then calculated a *P* value to assess whether laser stimulation significantly modulated the brain state.

Samples sizes. To determine the sample sizes of the experimental groups for optogenetic and pharmacogenetic experiments, for each group we first performed pilot experiments with two or three mice. Given the strength of the effect and the variance across this group, we then predicted the number of animals required to reach sufficient statistical power. To determine the sample size (number of units) for optrode recordings, we first recorded from two animals. Given the success rate of finding identified units and the homogeneity of units in the initial data set, we set a target sample size of 20 units. For rabies-mediated, anterograde, or dual retrograde tracing experiments, the selection of the sample size was based on numbers reported in previous studies. Otherwise, no statistical methods were used to predetermine sample size.

32. Miyamichi, K. *et al.* Dissecting local circuits: parvalbumin interneurons underlie broad feedback control of olfactory bulb output. *Neuron* **80**, 1232–1245 (2013).
33. Anikeeva, P. *et al.* Optetrode: a multichannel readout for optogenetic control in freely moving mice. *Nature Neurosci.* **15**, 163–170 (2012).
34. Schmitzer-Torbert, N., Jackson, J., Henze, D., Harris, K. & Redish, A. D. Quantitative measures of cluster quality for use in extracellular recordings. *Neuroscience* **131**, 1–11 (2005).

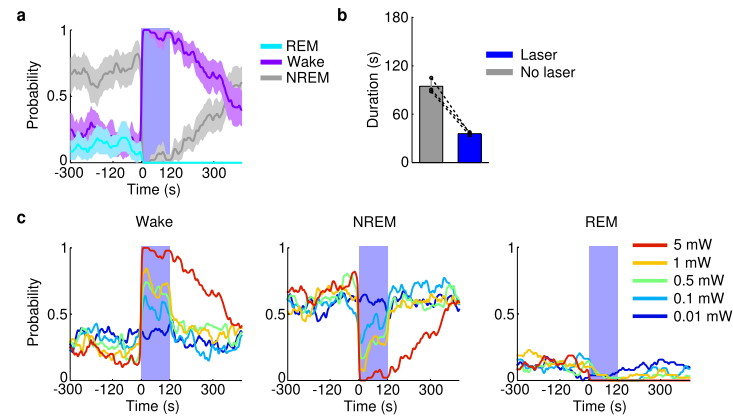


Extended Data Figure 1 | Comparison of spontaneous and laser-induced REM sleep. **a**, Mean EEG spectrogram and EMG amplitude before and after laser onset (averaged across all trials with laser onset falling on NREM sleep). **b**, Mean EEG spectrogram and EMG amplitude before and after spontaneous REM onset outside laser stimulation periods (only REM episodes with duration >70 s were included). **c**, Comparison of EEG power spectra during spontaneous (grey) and laser-induced (blue) REM sleep. Blue shading, s.e.m. for laser-induced REM sleep.



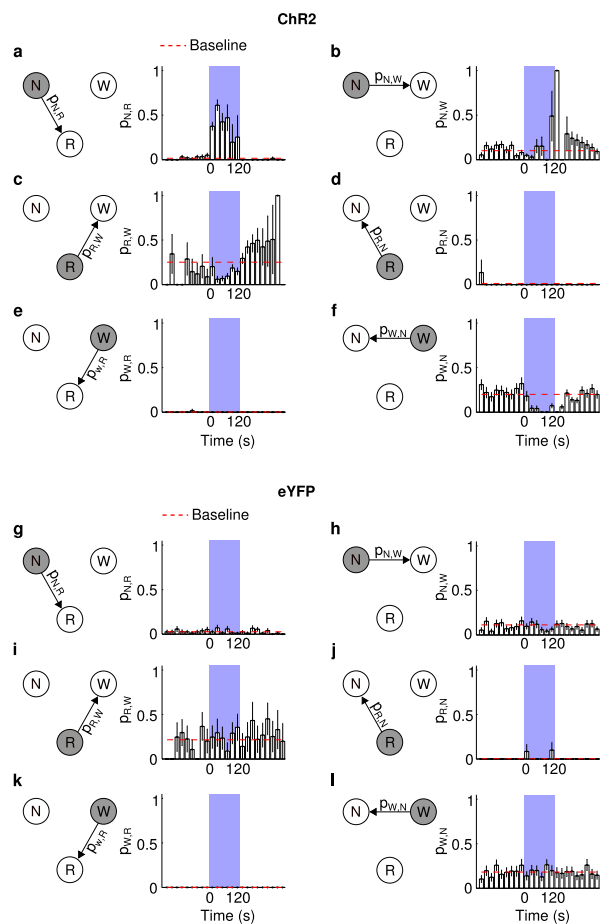
Extended Data Figure 2 | Effect of laser stimulation on brain states in eYFP control mice. **a**, Brain states in all trials from five mice aligned by time of laser stimulation (top) and probability of wake, NREM, or REM states before, during, and after laser stimulation (bottom). Shading, 95% CI. Blue bar, period

of laser stimulation. Laser stimulation caused no significant change in the probability of any brain state ($P > 0.34$, bootstrap). **b**, Similar to **a**, for trials in which laser onset fell on NREM sleep. **c**, Trials in which laser onset fell on wakefulness.

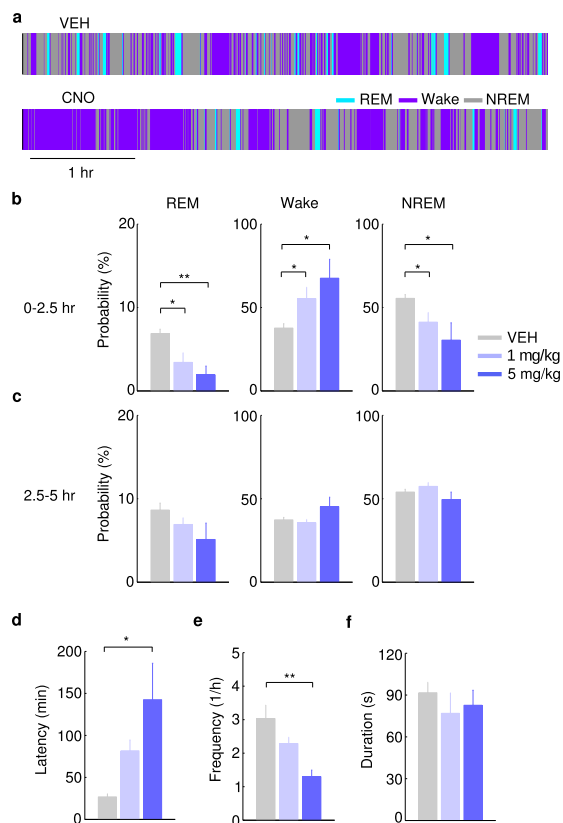


Extended Data Figure 3 | Optogenetic activation of vM glutamatergic neurons induces wakefulness. **a**, Probability of wake, NREM, or REM states before, during, and after laser stimulation (20 Hz, 120 s) in VGLUT2-Cre mice injected with AAV expressing ChR2-eYFP into the vM ($n = 3$ mice). Shading, 95% CI. Blue bar, period of laser stimulation. Laser stimulation caused a significant increase in wakefulness ($P < 0.001$, bootstrap) and decrease in

NREM sleep ($P < 0.001$). **b**, Mean durations of REM sleep episodes with and without laser stimulation. Each pair of dots represents data from one mouse. Laser stimulation shortened the duration of REM sleep episodes ($n = 3$ mice, $P = 0.008$, paired t -test). Error bar, s.d. **c**, Probability of wake (left), NREM (middle), and REM (right) states before, during, and after laser stimulation (20 Hz, 120 s) at different laser powers (colour coded).



Extended Data Figure 4 | Effect of laser stimulation of vM GABAergic neurons on the transition probability between each pair of brain states in Chr2 and eYFP control mice. a–f, Chr2 control mice. a, Probability of NREM (N) to REM (R) state transition within each 20 s period in Chr2 mice. Blue shading, period of laser stimulation (20 Hz, 120 s). Error bar, s.d. (bootstrap, $n = 6$ mice). Probability of baseline transition (red dashed line) was computed after excluding the laser stimulation period. The probability during laser stimulation was significantly higher than the baseline ($P = 0.02$, bootstrap). b, Similar to a, for NREM to wake (W) transition. The probability during laser stimulation was not significantly different from baseline ($P = 0.24$). c, REM to wake, probability during laser stimulation was significantly lower than baseline ($P < 0.001$), consistent with the effect of vM GABAergic neurons on prolonging REM duration (Fig. 2). d, REM to NREM, which rarely occurs in rodents. Laser stimulation caused no significant effect ($P > 0.99$). e, Wake to REM, which rarely occurs in normal mice. Laser stimulation had no significant effect ($P > 0.99$). f, Wake to NREM. Laser stimulation caused a significant reduction in the transition probability ($P < 0.001$), indicating that during wakefulness vM GABAergic neuron activity has a wake-maintenance effect. g–l, Similar to a–f, for eYFP control mice. Laser stimulation had no significant effect on any transition probability ($P > 0.05$).

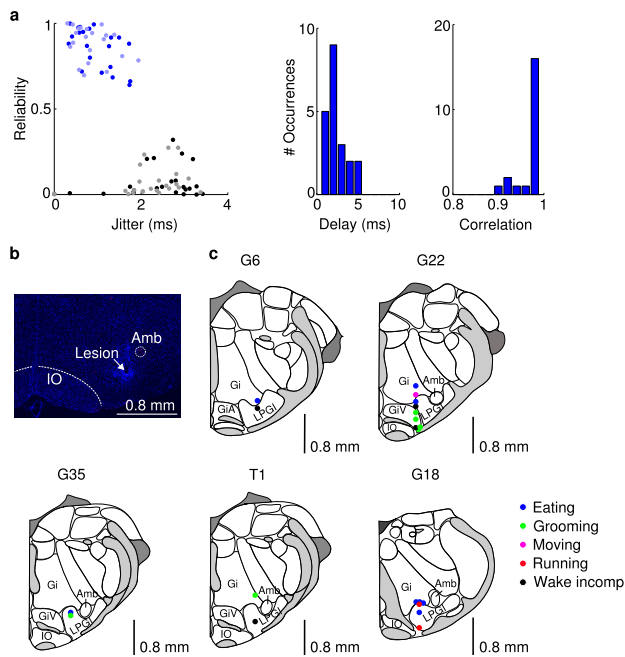


Extended Data Figure 5 | Pharmacogenetic inactivation of vM GABAergic neurons reduces REM sleep.

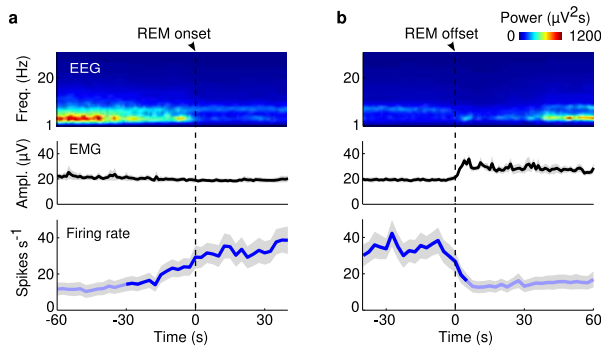
a, Brain states in a control (vehicle injection) and a CNO session from an example mouse. The recording session started 20 min after vehicle or CNO injection. **b**, Probability of each brain state during the first 2.5 h of the recording session, after injection of vehicle (grey) or two dosages of CNO (different shades of blue). Error bar, s.e.m. ($n = 6$ mice).

* $P < 0.05$; ** $P < 0.01$; one-way analysis of variance with post hoc Dunnett's test. **c**, Similar to **b**, but during the second half of the recording session (2.5–5 h).

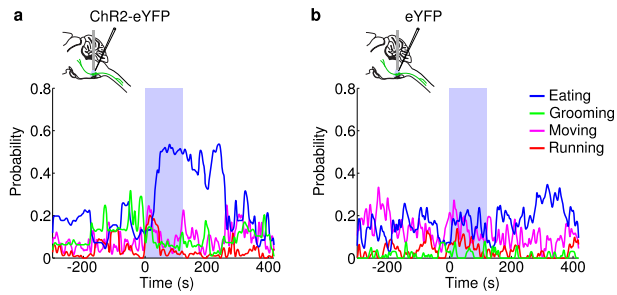
There was no significant difference between control and CNO at any dosage ($P > 0.12$). **d**, Latency of first REM sleep episode (from the beginning of each recording session). **e**, Frequency of REM episodes during the first 2.5 h of the recording session. **f**, Duration of REM episodes during the first 2.5 h of the session. The reduction of REM sleep caused by pharmacogenetic inactivation of vM GABAergic neurons appears to be due to the reduction of frequency rather than duration of REM episodes.



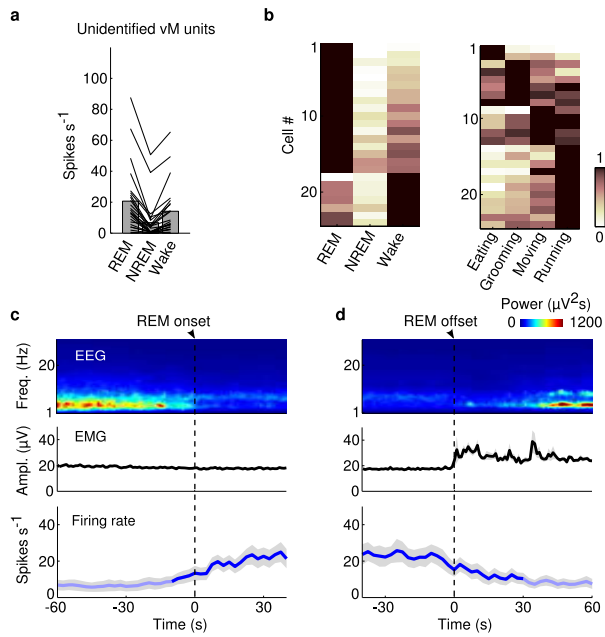
Extended Data Figure 6 | Optogenetic identification of vM GABAergic neurons. **a**, Left, reliability and temporal jitter of laser-evoked spikes in identified (blue) and unidentified (grey) units recorded in the vM. Note that the identified and unidentified units form two distinct clusters, with high reliability and low jitter for identified units. Dark/light symbols, data during 30/15 Hz laser stimulation. Middle, distribution of delays of laser-evoked spiking for 21 identified GABAergic neurons. Delay is defined as timing of the first spike after each laser pulse. Right, distribution of correlation coefficient between laser-evoked and spontaneous spike waveforms for all 21 identified vM GABAergic neurons. **b**, Fluorescence image of a coronal section showing the position of an electrolytic lesion at the end of the optrode tract (arrow). Blue, DAPI staining. **c**, Positions of the 21 identified vM GABAergic neurons from 5 mice. Each dot indicates one neuron. All 20 neurons with recording periods encompassing all three brain states showed maximal firing rates during REM sleep. The wakeful behaviour during which the neuron showed the maximal firing rate is colour-coded. Black, neurons for which the recording period did not include all wakeful behaviours. Schemes of brain sections adapted from Allen Mouse Brain Atlas (Website: © 2015 Allen Institute for Brain Science. Allen Mouse Brain Atlas [Internet]. Available from: <http://mouse.brain-map.org>).



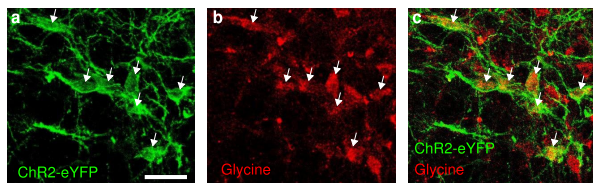
Extended Data Figure 7 | Activity of vM GABAergic neurons at REM sleep onset and offset. **a**, Mean firing rates of vM GABAergic neurons at REM onset. Shading, s.e.m. Dark blue, period in which firing rate was significantly higher than baseline ($P < 0.05$, Wilcoxon signed-rank test); baseline was defined as the average firing rate during the 10 s intervals 60 s before and after REM sleep. **b**, Mean firing rates at REM offset.



Extended Data Figure 8 | Effect of laser stimulation of vM neurons on several wakeful behaviours. **a**, Probability of moving, running, eating, and grooming before, during, and after laser stimulation (20 Hz, 120 s) in GAD2-Cre mice injected with AAV expressing ChR2-eYFP into the vM ($n = 8$ mice). Unclassified behaviours are not shown. Blue bar, period of laser stimulation. Laser stimulation caused a significant increase in eating ($P = 0.008$, Wilcoxon signed-rank test). **b**, Similar to **a**, but in control mice expressing eYFP ($n = 4$ mice).



Extended Data Figure 9 | Firing rates of unidentified vM neurons. **a**, Firing rates of unidentified units in the three brain states. Each line represents data from one neuron. Gray bar represents average over units ($n = 24$). **b**, Left, relative firing rates of the units across brain states. The firing rates of each unit were normalized by its maximum. Right, normalized firing rates across different wakeful behaviours. **c**, Mean firing rates of unidentified vM neurons at REM onset. Shading, s.e.m. Dark blue, period in which firing rate was significantly higher than baseline ($P < 0.05$, Wilcoxon signed-rank test). **d**, Mean firing rates of unidentified vM neurons at REM offset.



Extended Data Figure 10 | Co-expression of GABA and glycine in vM neurons. **a**, Fluorescence image of ChR2-eYFP expressing neurons in the vM of a GAD2-Cre mouse injected with Cre-inducible AAV. **b**, Immunohistochemical staining for glycine. **c**, Superposition of eYFP expression and glycine staining. In total, 94% (273/289) of eYFP-expressing cells were glycine positive ($n = 3$ mice).

Encoding of action by the Purkinje cells of the cerebellum

David J. Herzfeld¹, Yoshiko Kojima², Robijanto Soetedjo² & Reza Shadmehr¹

Execution of accurate eye movements depends critically on the cerebellum^{1–3}, suggesting that the major output neurons of the cerebellum, Purkinje cells, may predict motion of the eye. However, this encoding of action for rapid eye movements (saccades) has remained unclear: Purkinje cells show little consistent modulation with respect to saccade amplitude^{4,5} or direction⁴, and critically, their discharge lasts longer than the duration of a saccade^{6,7}. Here we analysed Purkinje-cell discharge in the oculomotor vermis of behaving rhesus monkeys (*Macaca mulatta*)^{8,9} and found neurons that increased or decreased their activity during saccades. We estimated the combined effect of these two populations via their projections to the caudal fastigial nucleus, and uncovered a simple-spike population response that precisely predicted the real-time motion of the eye. When we organized the Purkinje cells according to each cell's complex-spike directional tuning, the simple-spike population response predicted both the real-time speed and direction of saccade multiplicatively via a gain field. This suggests that the cerebellum predicts the real-time motion of the eye during saccades via the combined inputs of Purkinje cells onto individual nucleus neurons. A gain-field encoding of simple spikes emerges if the Purkinje cells that project onto a nucleus neuron are not selected at random but share a common complex-spike property.

Previous studies have focused on bursting activity of Purkinje cells during saccades^{6,10,11} and found no consistent modulation with saccade amplitude^{4,5}, speed^{5–7} or direction⁴. A recent simulation¹² suggested that Purkinje cells that pause during saccades may be important in understanding the responses observed in the deep cerebellar nucleus neurons. The main question that we wished to address was how the Purkinje cells encode the real-time motion of the eye.

We analysed simple-spike activity of 72 Purkinje cells in the oculomotor vermis (OMV, cerebellar lobules VI and VII) of five monkeys during saccades. The population included cells that exhibited increased activity (bursting; $n = 39$, Fig. 1a) or decreased activity (pausing; $n = 33$, Fig. 1b). Consistent with previous reports^{5,6,10}, most neurons were poorly modulated by saccade amplitude (Fig. 1c and Extended Data Fig. 1); however, the mean firing rate of burst cells (but not pause cells) increased significantly with saccade peak speed (Fig. 1d, $P < 10^{-10}$). Previous work had demonstrated that the population response encoded additional saccade-related information that was not reliably present in the responses of individual neurons^{6,13,14}. To examine the population response, we measured change in firing rates (from baseline) for the bursting and pausing cells during slow (400° s^{-1}) and fast (650° s^{-1}) saccades (Fig. 1e), pooled across all directions. The onset of change in firing rates in both populations generally led saccade onset by more than 50 ms. The termination of activity was also significantly later than the saccade: a 650° s^{-1} saccade was 38 ± 1.2 ms in duration (mean \pm s.e.m.), whereas activity of burst and pause cells persisted for more than 100 ms. Given that the cerebellum is thought to have a critical role in termination of ipsiversive saccades^{15,16}, it is unlikely that separate populations of burst or pause Purkinje cells control the motion of the eye, since their activity persists for much longer than the saccade.

Purkinje cells project to the caudal fastigial nucleus (cFN), where about 50 Purkinje cells converge onto a cFN neuron¹⁷. For each Purkinje cell we computed the probability of a simple spike in 1-ms time bins during saccades of a given peak speed, averaged across all

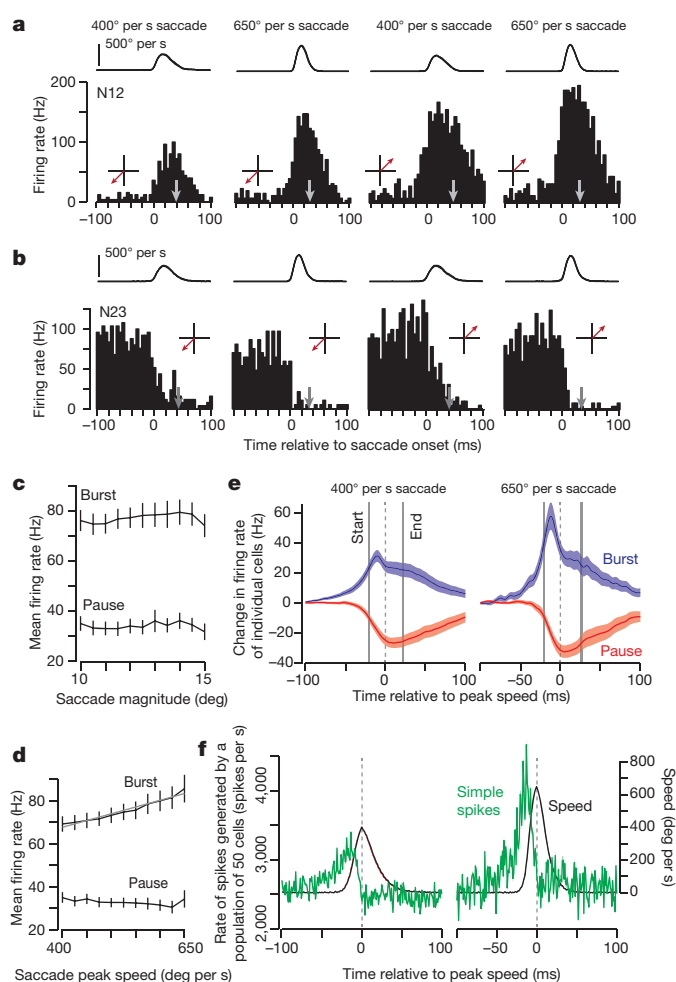


Figure 1 | Population of burst and pause Purkinje cells together predict eye speed in real time. **a, b**, Perisaccade histograms for a bursting (**a**) and pausing (**b**) Purkinje cell during saccades of various speeds and directions (red arrow). The trace on the top row is saccade speed. The grey arrow indicates saccade end. **c, d**, Mean firing rates over the duration of saccade computed across all directions. Changes in speed produced an increase in the firing rate of the burst cells but not the pause cells. **e**, Change in firing rates (with respect to baseline) of the bursting and pausing Purkinje cells for two saccade speeds. Grey bars are onset and termination of the saccade (width is s.e.m.). **f**, The total rate of simple spikes produced by a random selection of 50 Purkinje cells.

¹Department of Biomedical Engineering, Laboratory for Computational Motor Control, Johns Hopkins University School of Medicine, Baltimore, Maryland 21205, USA. ²Department of Physiology and Biophysics, Washington National Primate Center, University of Washington, Seattle, Washington 98195, USA.

directions. We then chose 50 Purkinje cells at random and computed the total number of simple spikes generated by the population at each millisecond, resulting in an estimate of the rate of presynaptic spikes converging onto a cFN cell. The results (Fig. 1f) revealed a real-time encoding of the speed of the eye: the peak of the activity preceded peak speed, increased in magnitude when speed increased, and returned to baseline just before saccade termination (R^2 at the optimal delay, 400° s^{-1} : $R^2 = 0.52$, $P < 10^{-22}$; 650° s^{-1} : $R^2 = 0.62$, $P < 10^{-43}$). It appeared that the simple spikes of the pause and burst cells combined together to predict motion of the eye.

Let us hypothesize that the Purkinje cells that project to a nucleus neuron are not selected randomly, but are organized by their inputs from the inferior olive¹⁸. That is, suppose that the olive projections divide the Purkinje cells into clusters where each cluster of Purkinje cells projects onto a single nucleus neuron. The input from the olive produces complex spikes in the Purkinje cells. We found that if we organized the simple spikes of the Purkinje cells based on each cell's complex-spike properties, additional features of the population activity were unmasked.

We measured complex-spike properties of each Purkinje cell by inducing a post-saccadic error through displacement of the target during the saccade, and then measured the probability of complex spikes as a function of the direction of this error (Fig. 2 and Supplementary Information section 2). For each Purkinje cell, the direction of error that produced the largest probability of complex spikes during the 50–200-ms post-saccade period was labelled as CS-on, and the opposite direction was labelled as CS-off (Extended Data Fig. 2). We then made the assumption that the Purkinje cells that projected onto a nucleus neuron all had the same CS-on direction (Fig. 3a). Under this assumption, we computed the rate of presynaptic simple spikes that a nucleus neuron would receive from the cluster of Purkinje cells (Supplementary Information section 3). We did this by convolving each Purkinje cell simple-spike train with a 2.5 ms standard-deviation normalized Gaussian, approximating the temporal characteristics of the inhibition produced in the nucleus neuron due to a simple spike in the Purkinje cell^{17,19}.

Figure 3b shows the change in population response from the baseline level when a saccade was made in the same direction as CS-off. The response rose above baseline before saccade onset, peaked before peak speed, and then returned to near baseline. The peak response scaled robustly with saccade amplitude (Fig. 3c, $R^2 = 0.93$, $P < 10^{-5}$). We observed a strong correspondence between the real-time population response and the real-time speed (Fig. 3d, lower plot, and Extended Data Fig. 3). The population response preceded eye speed by an average of 21.2 ± 0.4 ms (correlation analysis in the CS-off direction, mean \pm s.e.m.). Peak population response precisely predicted peak speed (Fig. 3e, $R^2 = 0.98$, $P < 10^{-7}$).

We took advantage of natural variability in saccades to test further the relationship between the population response and speed. We sorted all 10° saccades (direction CS-off) according to peak speed (Fig. 3f) and found that despite the constant amplitude, the population response precisely predicts the actual peak speed of the eye (Fig. 3g, $R^2 = 0.96$, $P < 10^{-7}$). Therefore, when the simple spikes were organized according to each cell's complex-spike directional preference, the population response for saccades of constant direction predicted nearly all of the variability in saccade peak speed.

No previous work, to our knowledge, had revealed how direction of a saccade is encoded in the activity of Purkinje cells. For the burst and pause cells, the mean and peak firing rates during the saccade did not vary as a function of direction (Extended Data Fig. 4). However, organizing the population response according to each Purkinje cell's complex-spike tuning preference revealed a clear encoding of direction: the peak response was greater if the saccades were in the same direction as CS-off as compared to CS-on (Fig. 4a, t -test, $P < 10^{-16}$). Indeed, the peak population response rose linearly as a function of peak speed in both directions, but with a larger gain when the saccade direction was congruent with CS-off (Fig. 4b,

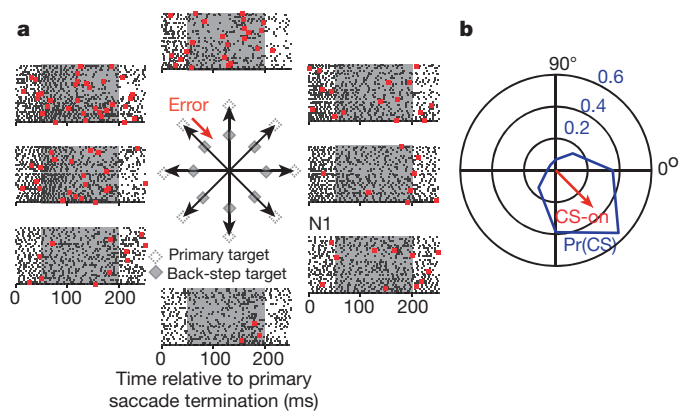


Figure 2 | Determination of complex-spike properties of Purkinje cells.

a, Response of a Purkinje cell during the 250-ms period after completion of a saccade (simple spikes are grey; complex spikes are red). CS-on was determined via a back-step paradigm in which the target was jumped (unfilled target to filled target) during saccade execution. Black arrow indicates saccade vector, red arrow indicates error vector. We computed the probability of complex spikes in the 50–200-ms period after saccade termination. **b**, The probability of a complex spike (Pr(CS)) as a function of the direction of the error vector. For this neuron, the highest probability (CS-on) occurred when the error vector was in direction -45° . The direction of CS-off for this cell was 135° .

repeated measures analysis of variance (ANOVA) with main effects of peak speed, $P < 10^{-15}$; CS-direction $P < 10^{-7}$; and a speed by CS-direction interaction, $P < 10^{-15}$).

To examine the effects of saccade direction more closely, we plotted the population response across saccade directions with respect to CS-on (Fig. 4c). We found that the population response was highest for saccades made in the CS-off direction, with an encoding of direction that resembled a cosine function (Fig. 4d). Therefore, the combined activity of burst and pause cells, but not the activity of either population individually (Extended Data Fig. 4), aligned to CS-off, produced a population response that exhibited gain-field encoding: the magnitude of the population response increased linearly with speed, and was cosine-tuned in direction, with a multiplicative interaction between speed and direction. The rate of simple spikes converging onto cFN, represented by $s(t)$, predicted in real-time motion of the eye (Supplementary Information section 4 and Extended Data Fig. 5) is:

$$s(t) = |\dot{\mathbf{x}}(t + \Delta)|g(\theta, \theta_{\text{CS}}) + c \quad (1)$$

$$g(\theta, \theta_{\text{CS}}) = a \cos(\theta - \theta_{\text{CS}}) + b$$

In equation (1), $|\dot{\mathbf{x}}(t + \Delta)|$ represents the magnitude of the eye velocity vector (the time derivative of eye position, \mathbf{x}) at time $t + \Delta$ (where $\Delta = 19$ ms), b and c are baseline offsets, a is a scaling factor, θ is saccade direction, and θ_{CS} is direction of CS-off for that cluster of Purkinje cells. The resulting gain-field encoding of eye motion is depicted in Fig. 4e.

We next addressed the question of how the activity of individual cells produced this directional encoding in the population response. The main contributors were the pause cells, which started their pause approximately 10 ms earlier when the saccade was in the CS-on direction (Fig. 4f), a change that was independent of saccade speed (Extended Data Fig. 6). This subtle shift in the timing of spikes produced an increase of the population response when saccade direction changed from CS-on to CS-off (Fig. 4a).

We found that the anatomical distribution of Purkinje cells, as labelled by their CS-off direction, was not random, but lateralized⁹ (Extended Data Fig. 7), confirming previous anatomical studies suggesting that olivary projections are contralateral^{20,21}. Purkinje cells with rightward CS-off were more likely to be on the right side of the cerebellum (t -test, $P < 10^{-4}$). This indicates that saccades made in the same direction as CS-off were typically ipsiversive, whereas saccades

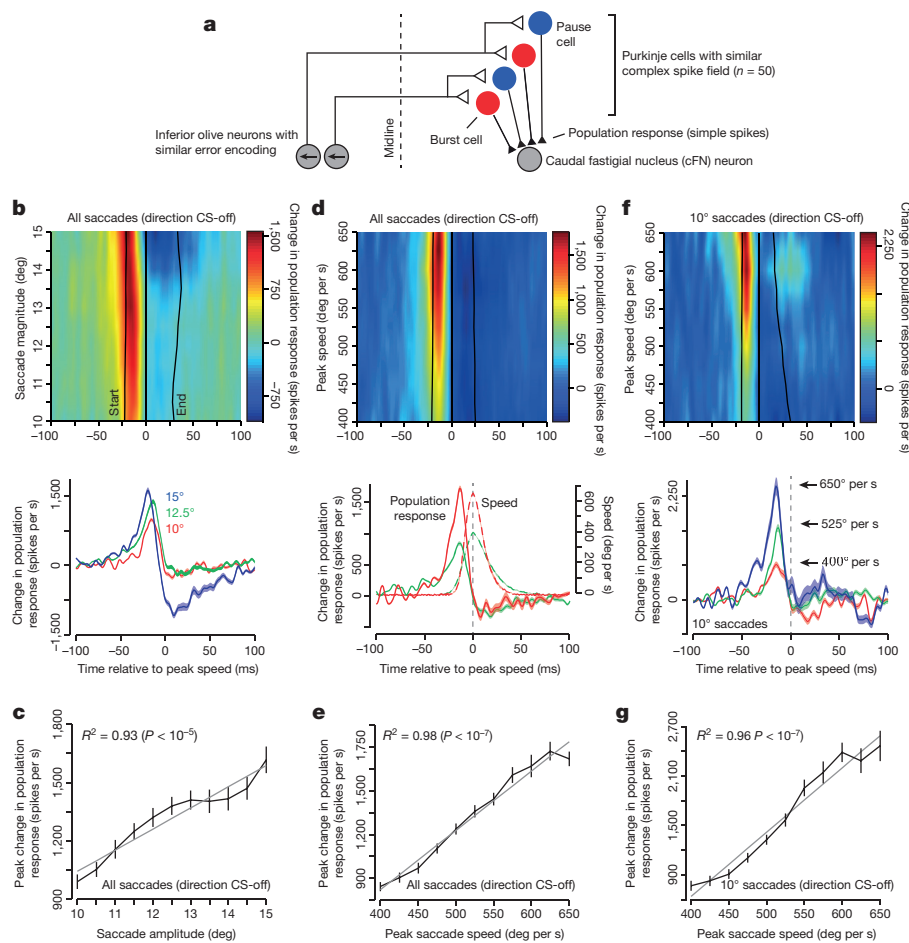


Figure 3 | A cluster of Purkinje cells, organized by their complex spikes, produced a population response that predicted in real time the motion of the eye. **a**, Hypothesized organization of the oculomotor vermis. To compute a population response, we measured the simple spikes of each Purkinje cell as a function of saccade direction with respect to the CS-on direction of that cell. For the Purkinje cells shown here, the CS-on is an error vector to the left (arrow). **b**, Change in population response (with respect to baseline) as a function of saccade amplitude in 0.5° bins, for saccades in the CS-off direction. Data in the amplitude axis were smoothed by a first-order Savitzky–Golay filter with a width of three bins⁶. Bottom plot shows the population response for three representative amplitudes. **c**, Peak population response increased linearly with saccade amplitude. *P* values indicate significant linear correlation. **d**, Population response as a function of saccade peak speed. Bottom plot shows representative responses with their corresponding speed traces. **e**, Peak population response increased linearly with saccade peak speed. **f**, Population response for 10° saccades ($\pm 1^\circ$), as a function of saccade peak speed. Bottom plot shows the population response for slow, medium and fast saccades of 10° amplitude. **g**, Peak population response increased linearly as a function of peak speed even for a fixed magnitude saccade. Error bars are s.e.m.

congruent with CS-on were contraversive. In contrast, pause and burst cells were uniformly distributed across the cerebellum ($P > 0.4$).

Our results rely critically on our hypothesis that Purkinje cells organize into clusters with roughly equal numbers of pause and burst

cells, all with a common complex-spike tuning preference (Fig. 3a). If, contrary to our hypothesis, pause and burst cells organized into separate clusters, the population response would not predict the real-time motion of the eye (Fig. 1e). Similarly, if each cluster was not

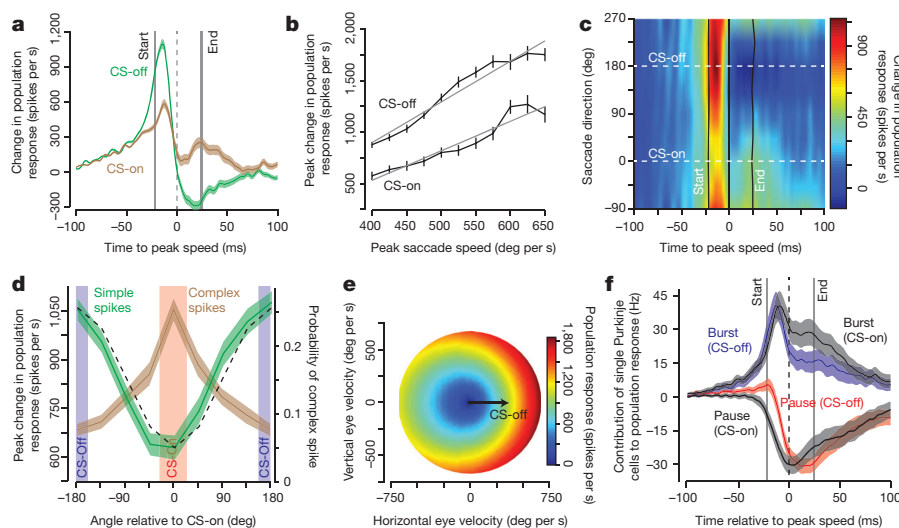


Figure 4 | Population response of Purkinje cells predicted saccade speed and direction in real time as a gain field. **a**, Population response for saccades in direction CS-on and CS-off. The population response is larger when the saccade is in the CS-off direction. **b**, Peak population response grew linearly with saccade speed, but had a higher gain for saccades in the CS-off direction. **c**, Real-time population response as a function of saccade direction relative to CS-on. Data smoothed as in Fig. 3b. **d**, Peak population response (labelled as

simple spikes) as a function of saccade direction with respect to CS-on. The brown curve shows probability of observing a complex spike as a function of the angle relative to each neuron's CS-on. The black curve indicates cosine fit of the peak population response. **e**, Gain-field encoding by a cluster of Purkinje cells whose CS-off direction is to the right (equation (1)). **f**, Contribution of single Purkinje cells to the population response. A change in direction coincides with a shift in timing of the pause cells. Error bars are s.e.m.

composed of roughly equal numbers of pause and burst cells, the population response could not predict the real-time speed of the eye (Extended Data Fig. 8 and Supplementary Information section 5). The fact that burst and pause cells were distributed uniformly across the recording locations, and not lateralized as we found with the complex-spike tuning properties, suggests that a cluster is composed of both burst and pause. Finally, if we ignored the complex-spike properties of the Purkinje cells, and made the typical assumption that simple spikes were sufficient to uncover the coordinate system of encoding motion, then the gain-field representation of speed and direction would disappear (Extended Data Fig. 9 and Supplementary Information section 6).

Organizing the Purkinje cell into clusters where all the cells shared a common complex-spike property resulted in simple spikes that encoded speed and direction in real time via a gain field.

Together, our results suggest three principles of cerebellar function during control of saccadic eye movements. First, the cerebellum predicts real-time motion not in the time course of individual Purkinje cell simple spikes, nor in the individual activities of the bursting or pausing populations, but in the combined activities of these two populations via the simple spikes that converge onto cells in the deep cerebellar nucleus. A similar population coding has been suggested during smooth pursuit²². Second, this population input to each nucleus neuron encodes direction and speed via a gain field. Because a similar encoding has been shown in the posterior parietal cortex during saccades²³, as well as in the motor cortex during reaching²⁴, our observation in the cerebellum suggests a common principle of encoding in disparate regions of the motor system. Finally, the gain-field encoding was present if we assumed a specific anatomical organization: a cluster of Purkinje cells that projected onto a single nucleus neuron was composed of approximately equal numbers of bursting and pausing Purkinje cells, all sharing a common complex-spike property. Because the complex spikes of a Purkinje cell are due to input from the inferior olive, the gain-field encoding predicts that the oculomotor vermis is organized into clusters of Purkinje cells that share similar climbing fibre projections from the inferior olive. This, in turn, suggests that motor memories are anatomically clustered in the cerebellum by the errors that were experienced during movements²⁵.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 24 February; accepted 8 September 2015.

- Xu-Wilson, M., Chen-Harris, H., Zee, D. S. & Shadmehr, R. Cerebellar contributions to adaptive control of saccades in humans. *J. Neurosci.* **29**, 12930–12939 (2009).
- Barash, S. *et al.* Saccadic dysmetria and adaptation after lesions of the cerebellar cortex. *J. Neurosci.* **19**, 10931–10939 (1999).
- Kojima, Y., Soetedjo, R. & Fuchs, A. F. Effects of GABA agonist and antagonist injections into the oculomotor vermis on horizontal saccades. *Brain Res.* **1366**, 93–100 (2010).
- Ohtsuka, K. & Noda, H. Discharge properties of Purkinje cells in the oculomotor vermis during visually guided saccades in the macaque monkey. *J. Neurophysiol.* **74**, 1828–1840 (1995).
- Helmchen, C. & Büttner, U. Saccade-related Purkinje cell activity in the oculomotor vermis during spontaneous eye movements in light and darkness. *Exp. Brain Res.* **103**, 198–208 (1995).

- Thier, P., Dicke, P. W., Haas, R. & Barash, S. Encoding of movement time by populations of cerebellar Purkinje cells. *Nature* **405**, 72–76 (2000).
- Kase, M., Miller, D. C. & Noda, H. Discharges of Purkinje cells and mossy fibres in the cerebellar vermis of the monkey during saccadic eye movements and fixation. *J. Physiol. (Lond.)* **300**, 539–555 (1980).
- Kojima, Y., Soetedjo, R. & Fuchs, A. F. Changes in simple spike activity of some Purkinje cells in the oculomotor vermis during saccade adaptation are appropriate to participate in motor learning. *J. Neurosci.* **30**, 3715–3727 (2010).
- Soetedjo, R., Kojima, Y. & Fuchs, A. F. Complex spike activity in the oculomotor vermis of the cerebellum: a vectorial error signal for saccade motor learning? *J. Neurophysiol.* **100**, 1949–1966 (2008).
- Catz, N., Dicke, P. W. & Thier, P. Cerebellar-dependent motor learning is based on pruning a Purkinje cell population response. *Proc. Natl Acad. Sci. USA* **105**, 7309–7314 (2008).
- Catz, N., Dicke, P. W. & Thier, P. Cerebellar complex spike firing is suitable to induce as well as to stabilize motor learning. *Curr. Biol.* **15**, 2179–2189 (2005).
- Gad, Y. P. & Anastasio, T. J. Simulating the shaping of the fastigial deep nuclear saccade command by cerebellar Purkinje cells. *Neural Netw.* **23**, 789–804 (2010).
- Dash, S., Dicke, P. W. & Thier, P. A vermal Purkinje cell simple spike population response encodes the changes in eye movement kinematics due to smooth pursuit adaptation. *Front. Syst. Neurosci.* **7**, 3 (2013).
- Pra, M., Dash, S., Catz, N., Dicke, P. W. & Thier, P. Characteristics of responses of Golgi cells and mossy fibers to eye saccades and saccadic adaptation recorded from the posterior vermis of the cerebellum. *J. Neurosci.* **29**, 250–262 (2009).
- Robinson, F. R., Straube, A. & Fuchs, A. F. Role of the caudal fastigial nucleus in saccade generation. II. Effects of muscimol inactivation. *J. Neurophysiol.* **70**, 1741–1758 (1993).
- Fuchs, A. F., Robinson, F. R. & Straube, A. Role of the caudal fastigial nucleus in saccade generation. I. Neuronal discharge pattern. *J. Neurophysiol.* **70**, 1723–1740 (1993).
- Person, A. L. & Raman, I. M. Purkinje neuron synchrony elicits time-locked spiking in the cerebellar nuclei. *Nature* **481**, 502–505 (2012).
- De Zeeuw, C. I. *et al.* Spatiotemporal firing patterns in the cerebellum. *Nature Rev. Neurosci.* **12**, 327–344 (2011).
- Telgkamp, P., Padgett, D. E., Ledoux, V. A., Woolley, C. S. & Raman, I. M. Maintenance of high-frequency transmission at Purkinje to cerebellar nuclear synapses by spillover from boutons with multiple release sites. *Neuron* **41**, 113–126 (2004).
- Yamada, J. & Noda, H. Afferent and efferent connections of the oculomotor cerebellar vermis in the macaque monkey. *J. Comp. Neurol.* **265**, 224–241 (1987).
- Kralj-Hans, I., Baizer, J. S., Swales, C. & Glickstein, M. Independent roles for the dorsal paraflocculus and vermal lobule VII of the cerebellum in visuomotor coordination. *Exp. Brain Res.* **177**, 209–222 (2006).
- Krauzlis, R. J. Population coding of movement dynamics by cerebellar Purkinje cells. *Neuroreport* **11**, 1045–1050 (2000).
- Andersen, R. A., Essick, G. K. & Siegel, R. M. Encoding of spatial location by posterior parietal neurons. *Science* **230**, 456–458 (1985).
- Paninski, L., Shoham, S., Fellows, M. R., Hatsopoulos, N. G. & Donoghue, J. P. Superlinear population encoding of dynamic hand trajectory in primary motor cortex. *J. Neurosci.* **24**, 8551–8561 (2004).
- Herzfeld, D. J., Vaswani, P. A., Marko, M. K. & Shadmehr, R. A memory of errors in sensorimotor learning. *Science* **345**, 1349–1353 (2014).

Supplementary Information is available in the online version of the paper.

Acknowledgements These data were collected in the laboratory of A. Fuchs. The authors are grateful to his generosity. The authors would like to thank S. du Lac for comments. The work was supported by NIH grants R01NS078311, R01EY019258, R01EY023277 and F31NS090860.

Author Contributions Y.K. and R.S. conceived, designed and performed all experiments. D.J.H. and R.Sh. formed the conceptual model. D.J.H. analysed the data and made all figures. R.Sh. and D.J.H. wrote the paper.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to D.J.H. (dherzfe1@jhmi.edu) or R.S. (shadmehr@jhu.edu).

METHODS

No statistical methods were used to predetermine sample size.

We analysed extracellular recordings from Purkinje cells of the oculomotor vermis in five rhesus monkeys (B, F, W, K, KO) while they made saccades to visual targets^{8,9} (Supplementary Information section 1). Each cell was well isolated for an average of 3,000 saccades. Briefly, a scleral search coil was surgically attached to the eye of each monkey, allowing measurement of eye position²⁶ while the animal's head was restrained. After surgery, the monkeys were trained to make saccades to targets of varying amplitudes and directions. We identified Purkinje cell activity in OMV by their saccade-related change in the simple-spike response, as well as the presence of complex spikes. Neurophysiological data were sampled at 50 kHz. All experiments were performed in accordance with the Guide for the Care and Use of Laboratory Animals (1997) and exceeded the minimal requirements recommended by the Institute of Laboratory Animal Resources and the Association for Assessment and Accreditation of Laboratory Animal Care International. All animal procedures were approved by the local committee at the University of Washington.

The CS-on direction for each cell was determined using the standard intra-saccadic step paradigm²⁷, in which the target was displaced during the initial saccade (Fig. 2). This error resulted in complex spikes during the period following the saccade, when the monkey observed the error. For every cell, we determined the CS-on direction as the error direction which elicited the largest probability of complex spikes during the 50–200 ms following the primary saccade. For $n = 39$ cells we were able to maintain excellent isolation of the Purkinje cell throughout the experiment, allowing us to perform automated identification of complex spikes on every trial. This allowed us to compute the probability of complex spikes as a function of error direction. For the remaining $n = 33$ cells, the CS-on direction was determined via analysis of the initial 50 trials for each direction of error⁸. CS-off was defined as CS-on + 180°.

We computed firing rates by determining the inverse of the time between two consecutive simple spikes²⁸ and then convolved the resulting time series with a normalized Gaussian kernel with a standard deviation of 2.5 ms, which is significantly shorter than conventional kernels⁶. This guarded against overestimating the duration of a population response. We calculated the mean firing rate during the

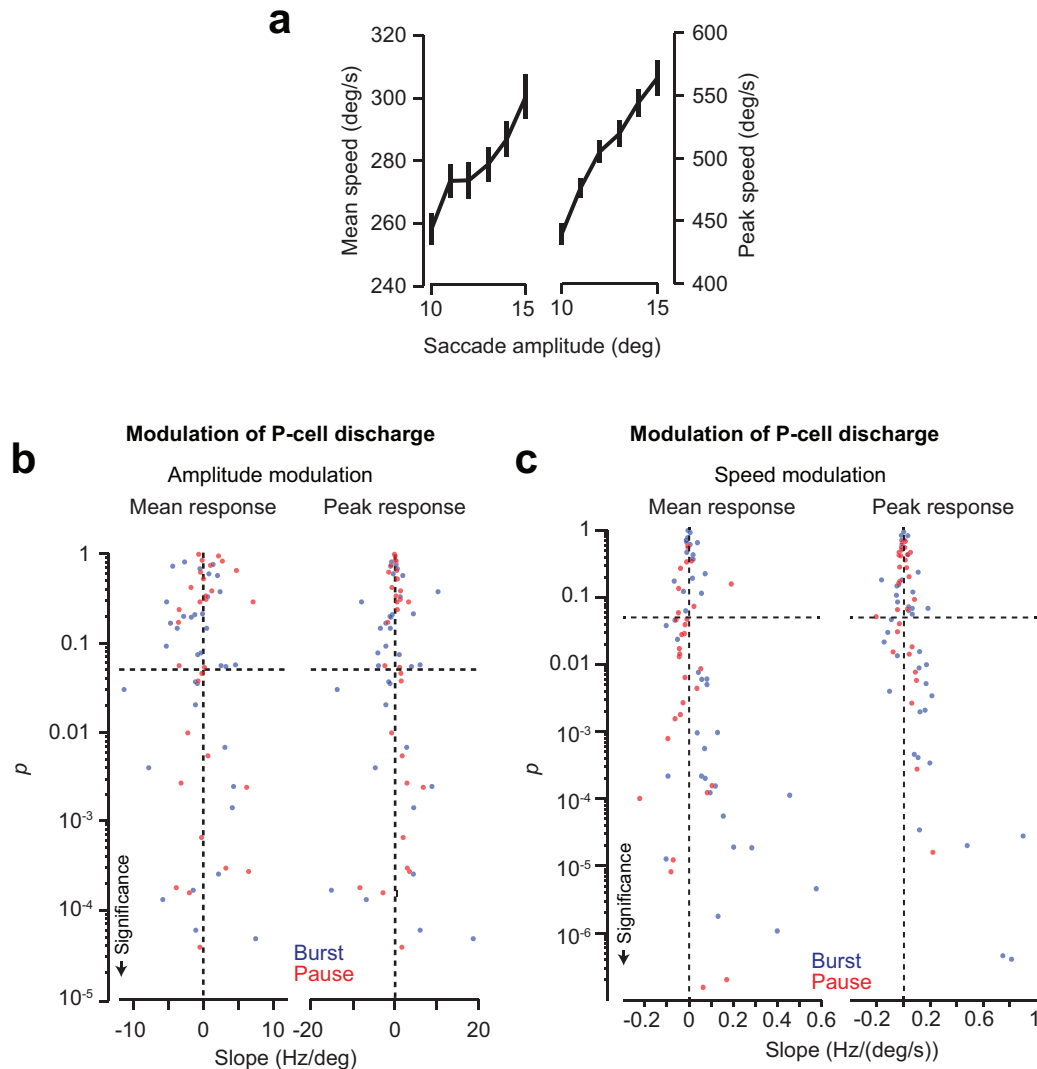
saccade by taking the average of the firing rate from the start to the end of the saccade. We determined the peak response of a Purkinje cell as the maximum firing during the saccade period.

The response of an individual Purkinje cell is quite variable⁸, with some neurons showing a combination of bursting and pausing activity in the time period near saccade onset. Therefore, to categorize neurons exclusively as pausing or bursting, we compared the mean firing rate of each cell from the period 200–100 ms before saccade onset of all recorded trials to a 150-ms period centred on saccade peak speed. Neurons that reduced their rate during this extended period were classified as 'pausing', whereas neurons that significantly exceeded this rate were classified as 'bursting'. We tested this categorization statistically via a paired t -test with a cutoff of $P = 0.05$. Only two neurons (both bursting) did not pass this statistical test. We included these two neurons in the bursting population.

To establish confidence intervals on the population responses, we performed a bootstrap analysis in which we randomly sampled 50 neurons from the available pool of 72 (with replacement), which simulated the approximate number of Purkinje cell inputs that project onto a nucleus neuron¹⁷. Error bars show mean \pm s.e.m. of 50 bootstrapped Purkinje cell populations. In cases where we show the responses of the bursting/pausing populations separately, we report the mean \pm s.e.m. across neurons in the respective population without bootstrapping.

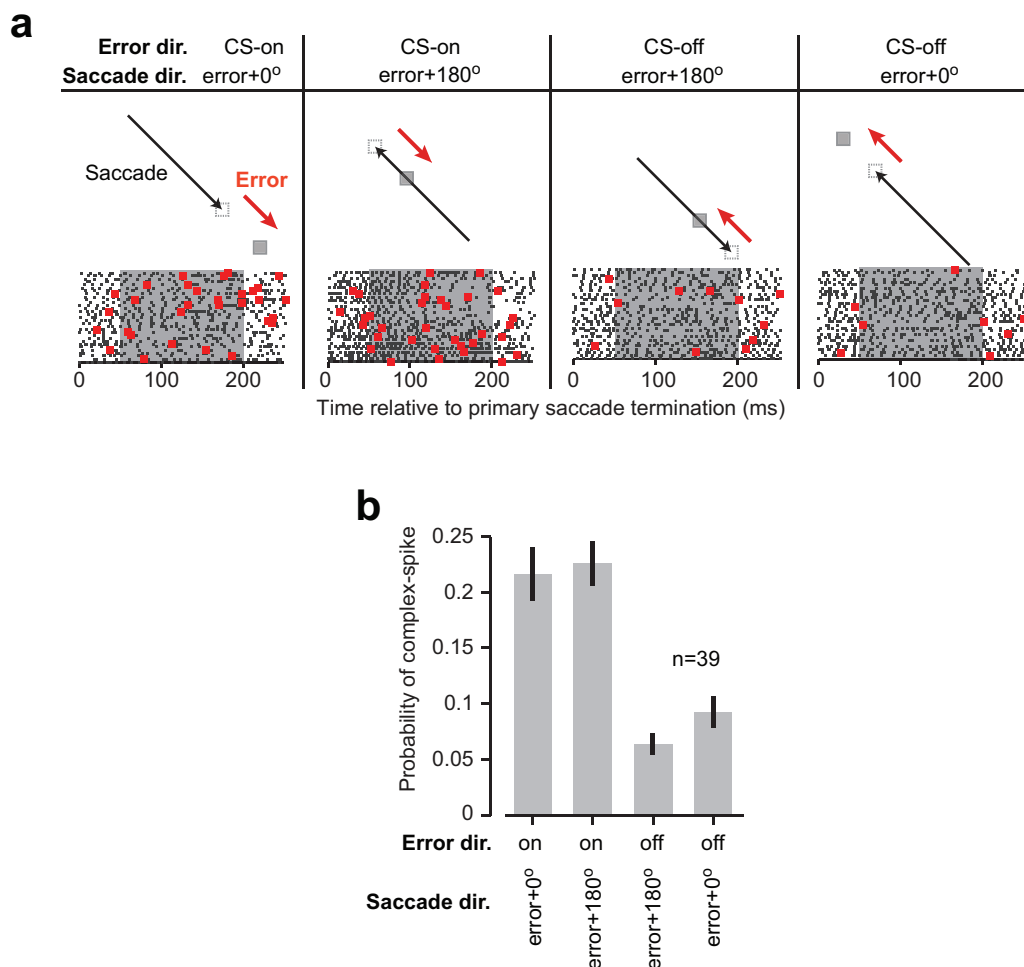
To compute the population response of a cluster of 50 Purkinje cells, we first convolved the simple spikes of each Purkinje cell with a kernel (normalized Gaussian of 2.5 ms s.d.), approximately representing the time course of post-synaptic inhibition induced by the simple-spike train¹⁷. We then computed the change from baseline for each cell, and finally the sum of changes across the population, resulting in a population response that had units of spikes s^{-1} , computed at each millisecond of time.

26. Fuchs, A. F. & Robinson, D. A. A method for measuring horizontal and vertical eye movement chronically in the monkey. *J. Appl. Physiol.* **21**, 1068–1070 (1966).
27. McLaughlin, S. C. Parametric adjustment in saccadic eye movements. *Percept. Psychophys.* **2**, 359–362 (1967).
28. Lisberger, S. G. & Pavelko, T. A. Vestibular signals carried by pathways subserving plasticity of the vestibulo-ocular reflex in monkeys. *J. Neurosci.* **6**, 346–354 (1986).



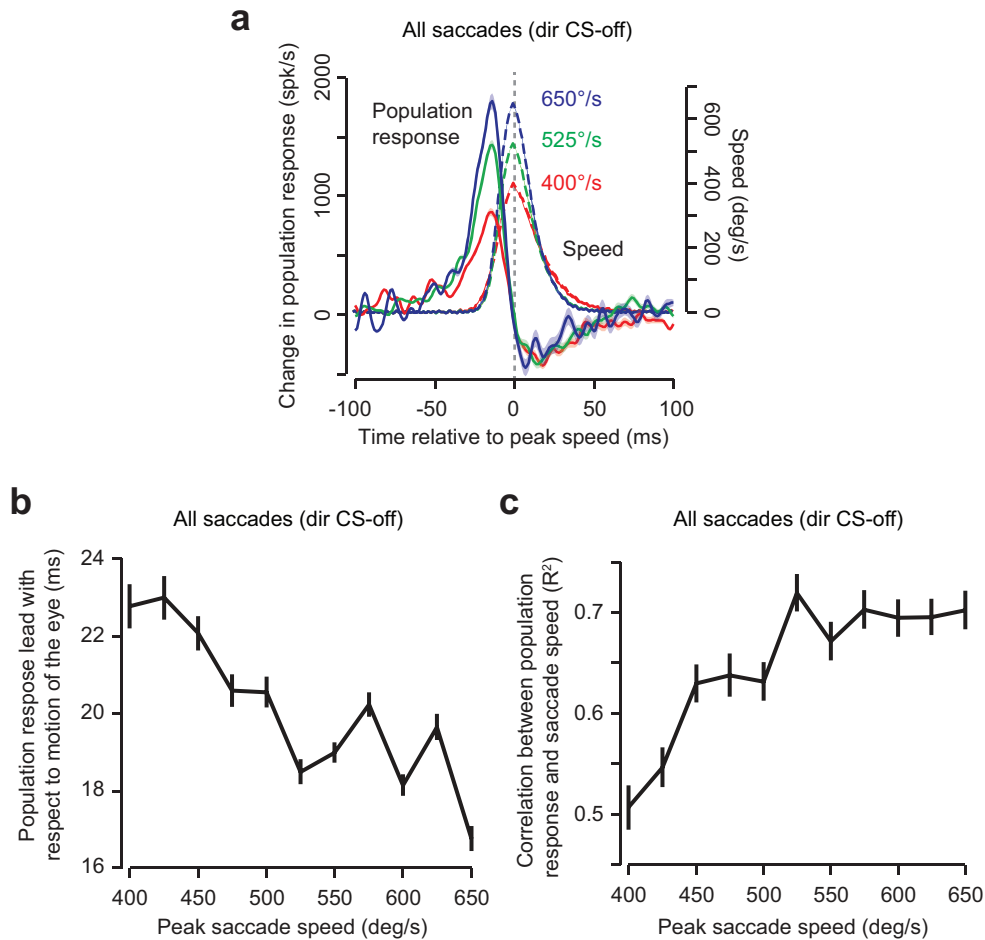
Extended Data Figure 1 | Firing rates of individual Purkinje cells as a function of saccade amplitude and peak speed. **a**, Increase in saccade amplitude produced robust increases in mean and peak saccade speed (mean: $R^2 = 0.86$, $P < 10^{-4}$; peak: $R^2 = 0.99$, $P < 10^{-9}$). Error bars indicate s.e.m. **b**, For each neuron, we correlated the average firing rate and the peak firing rate (computed over the saccade duration and averaged over all directions) with saccade amplitude. Some neurons increased their firing rates with increasing saccade amplitude (positive slope) and some neurons decreased

their responses (negative slope). However, mean and peak firing rates of a majority of neurons (47 of 72) were not significantly modulated with saccade amplitude. As a result, activity of neither the burst nor the pause cells showed a significant modulation with saccade amplitude (Fig. 1c, main manuscript). **c**, Most neurons (45 of 72) had a significant linear relationship between firing rates and peak saccade speed. In particular, mean and peak response of burst cells showed a significant increase with peak speed (Fig. 1d).



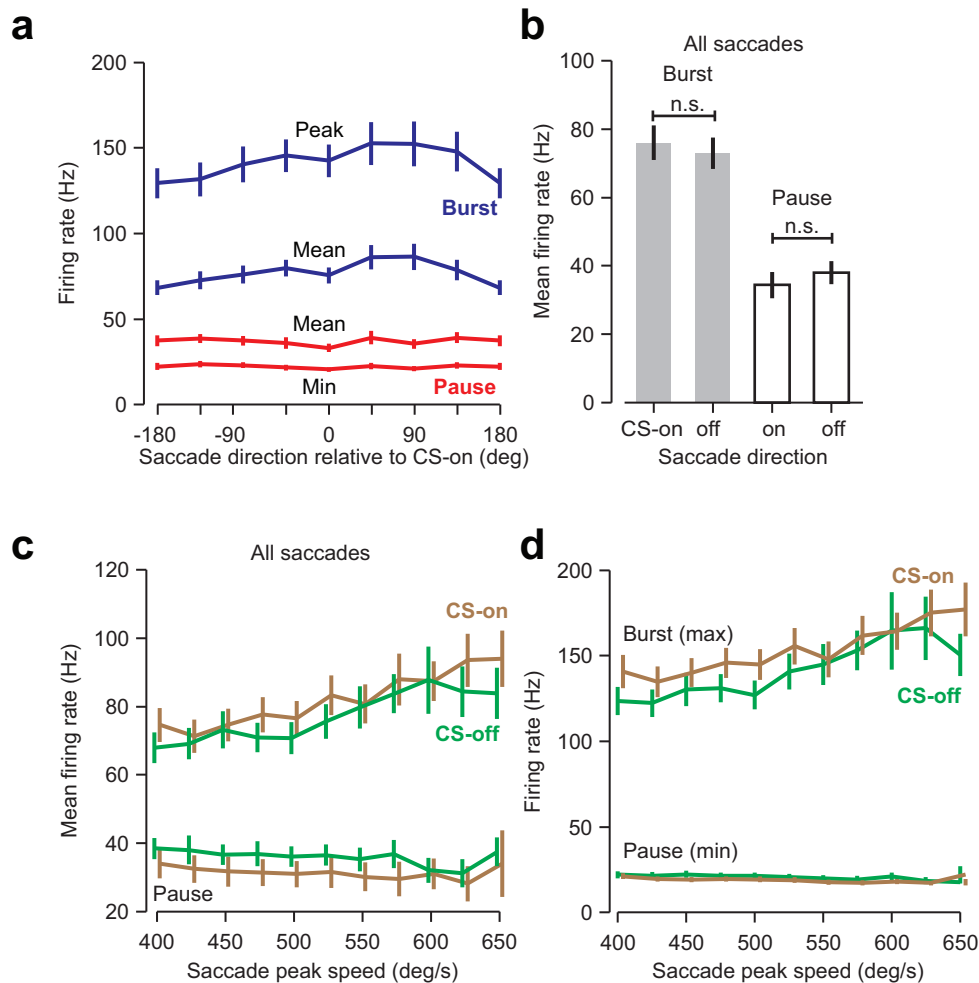
Extended Data Figure 2 | Complex spikes encode direction of error, not direction of saccade that preceded that error. **a**, The response of the same cell shown in Fig. 2 as a function of direction of saccade and direction of error. The probability of complex spikes is high when the direction of error is at -45° , despite the fact that saccade direction may be at -45° or $+135^\circ$. **b**, Population

statistics from $n = 39$ Purkinje cells in which the probability of complex spikes was quantified as a function of direction of the error vector and direction of the saccade that preceded that error. Probability of complex spikes depended on direction of error, not direction of saccade.



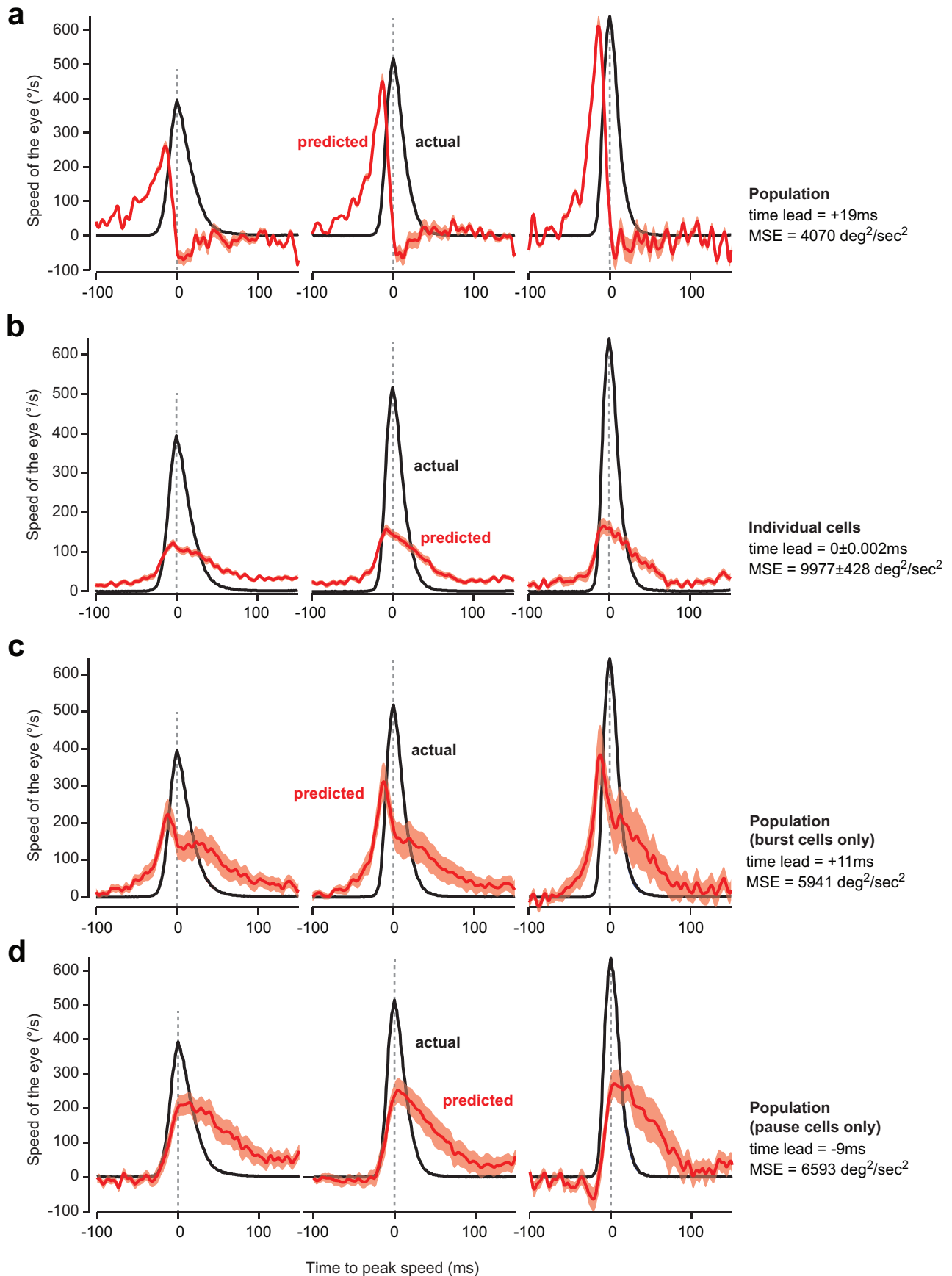
Extended Data Figure 3 | The simple-spike population response of Purkinje cells, organized by their complex-spike properties (Fig. 3a), correlated with motion of the eye in real time. a, Population response for saccades in direction CS-off for three different peak speeds. **b,** Temporal lead of the

population response with respect to saccade speed as computed by finding the temporal shift that maximized the cross-correlation. **c,** Correlation between the population response and the temporally shifted eye speed trace (measured as R^2). Error bars in all panels indicate s.e.m. across bootstrapped populations.



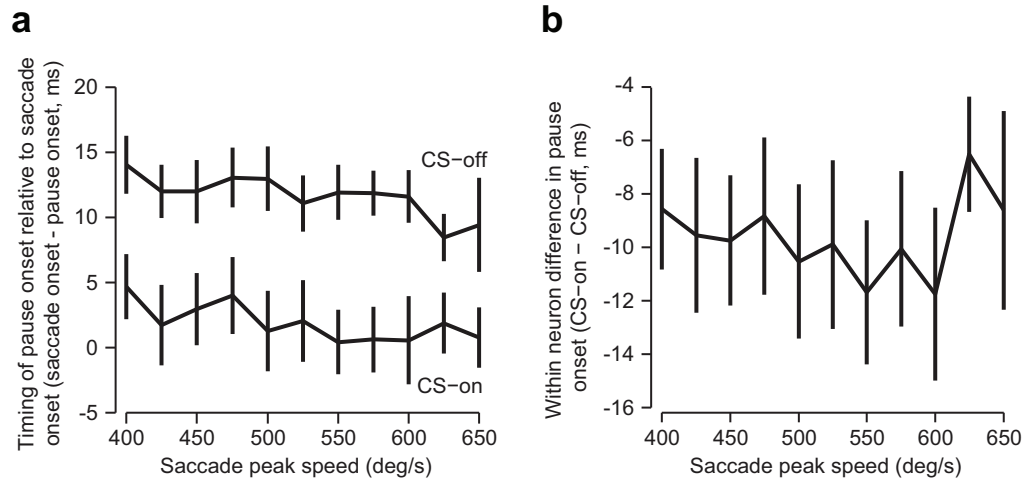
Extended Data Figure 4 | Mean and peak/trough firing rate of the burst and pause cells were poorly modulated by saccade direction. **a**, Maximum, minimum and mean firing rates averaged across burst or pause cells with respect to saccade direction, relative to CS-on direction of each cell. **b**, Mean firing rates of the burst and pause cells, as measured across all saccades, were not significantly different for saccades in the CS-on versus CS-off direction (burst $P > 0.10$, pause $P > 0.05$). **c**, Mean firing rates of the burst and pause cells as a function of saccade speed, for saccades in the CS-on versus CS-off direction. Saccade speed modulated the mean firing rates of the burst cells, but there were no significant interaction between saccade direction and speed ($P > 0.6$), nor a significant effect of saccade direction ($P > 0.7$). **d**, Peak (maximum) firing rates of the burst cells and the minimum firing rate of the

pause cells as a function of saccade speed, for saccades in the CS-off and CS-on directions. We asked whether the maximum response of the burst cells or the minimum response of the pause cells was significantly modulated by direction. Separate repeated measures ANOVAs showed that for the burst cells, peak activity increased as a function of saccade peak speed ($P < 0.001$), but this relationship was unaffected by saccade direction ($P > 0.4$). For the pause cells, the response was not affected by saccade speed ($P > 0.6$), and this relationship was not modulated by saccade direction ($P > 0.4$). We found that saccade direction did not significantly alter the encoding of peak speed in either the mean or minimum/maximum activity of Purkinje cells. Error bars in all panels represent s.e.m. across neurons.



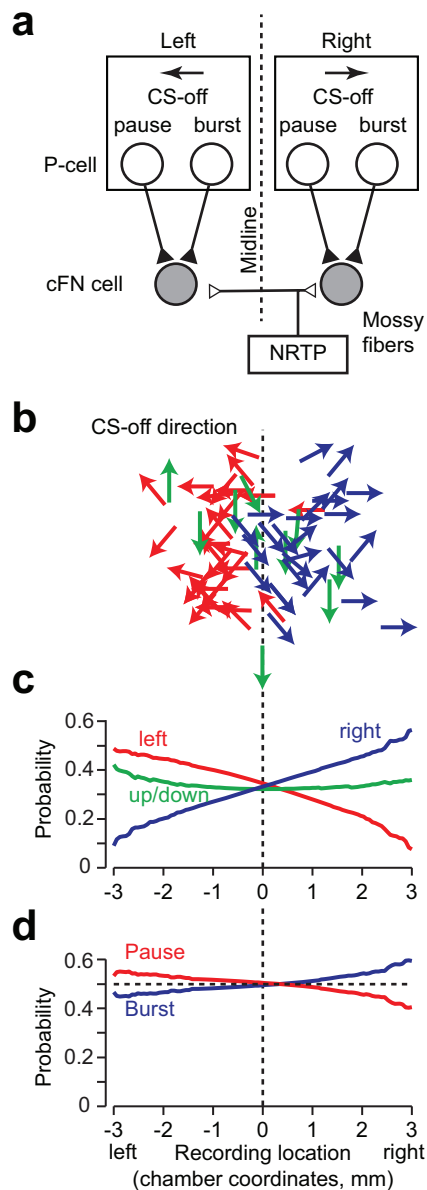
Extended Data Figure 5 | A population of Purkinje cells, organized by their complex spike properties, predicted the real-time speed of the eye better than activity of individual cells. **a**, We used equation (S2) (see Supplementary Information for details) and used the measured population response $s(t)$ of Purkinje cells to predict the real-time speed of the eye $|\hat{\mathbf{x}}(t + \Delta)|$. The plot shows the predicted speed for saccades of 400, 525 and 650° s⁻¹. The predicted speed led the actual speed by 19 ms. MSE is the mean squared error between the

predicted and actual eye trajectory at the optimal value of Δ . **b**, The result of fitting equation (S2) (see Supplementary Information) to the response of individual neurons. **c**, The result of fitting equation (S2) (Supplementary Information) to the discharge of a population composed exclusively of burst cells. **d**, The result of fitting equation (S2) (Supplementary Information) to the discharge of a population composed exclusively of pause cells.

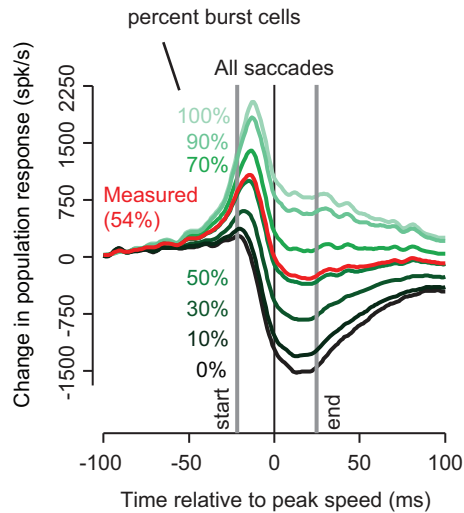


Extended Data Figure 6 | Change in saccade direction was associated with a change in the timing of the reduction of discharge in the pause cells (that is, pause onset) (see Fig. 4f). a, Timing of pause onset with respect to saccade onset for saccades of various speeds and directions. We computed the pause onset as the time when the neuron's response reached 20% of its minimum response. Positive numbers indicate that the pause onset occurred before

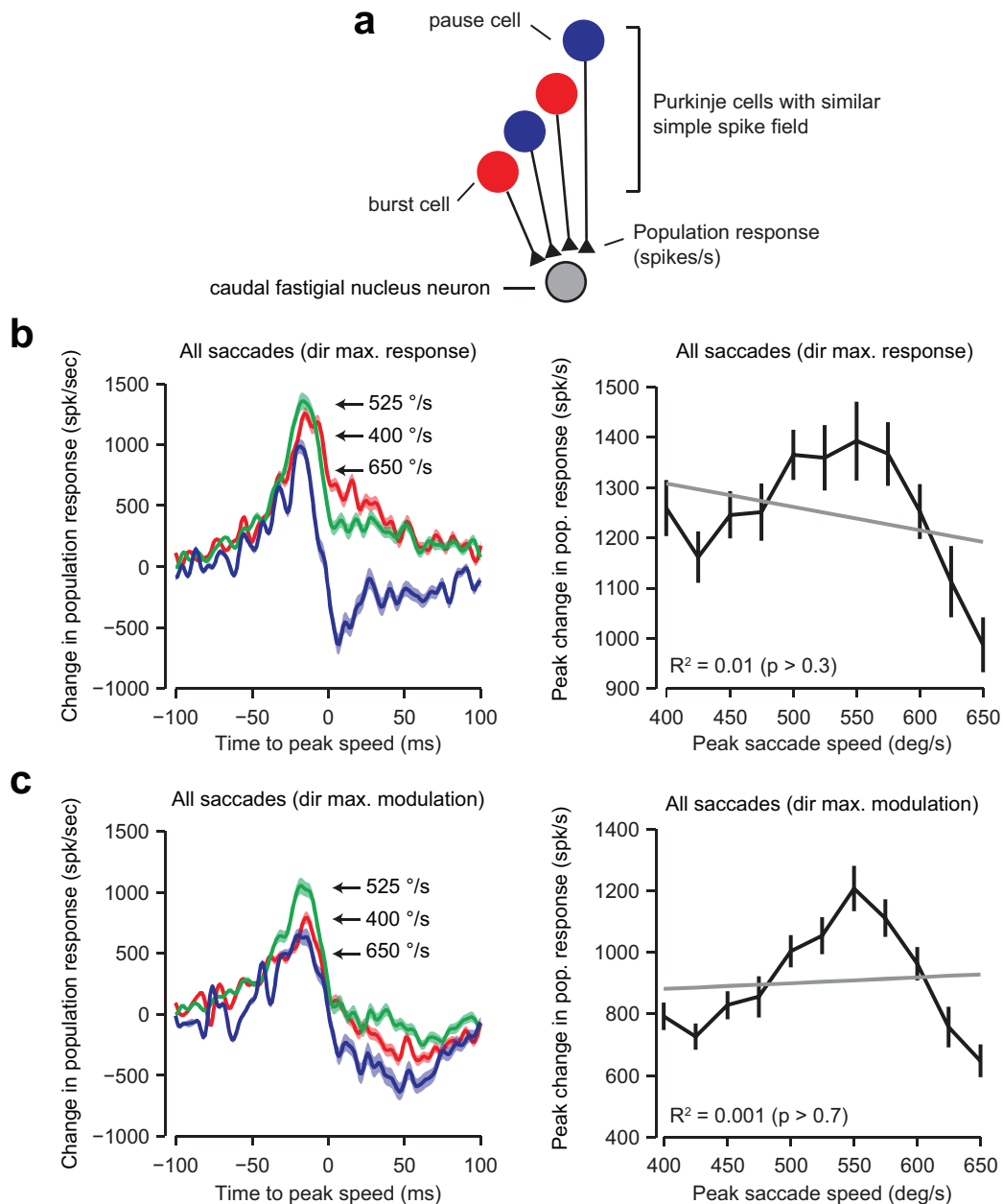
saccade onset. **b,** Within-neuron measure of pause onset for saccade in direction CS-on, minus onset from saccades in direction CS-off. Negative numbers indicate that the pause onset occurred earlier for saccades in the CS-on direction. Error bars in all panels indicate s.e.m. across neurons.



Extended Data Figure 7 | Complex-spike-dependent organization of the Purkinje cells. **a**, Hypothesized anatomical organization of the oculomotor vermis (OMV). Bursting and pausing Purkinje cells are organized into clusters, with the cells in each cluster sharing a common complex-spike direction. Neurons on the right side of the OMV project to right cFN neurons and have CS-off directions to the right. **b**, Distribution of the CS-off directions from recorded neurons in chamber coordinates. Vertical dotted line shows the line that best separates rightwards CS-off direction neurons (blue) from leftwards CS-off direction neurons (red). **c**, Probability of having a rightwards (blue), up/down (green), or leftwards CS-off direction as a function of chamber coordinates. Purkinje cells with CS-off to the left were more probable on the left side of the cerebellum. Purkinje cells with CS-off to the right were more probable on the right side of the cerebellum. **d**, Pause (red) and burst (blue) Purkinje cells were equally likely at all recorded locations.



Extended Data Figure 8 | The population response was sensitive to the fraction of pause and burst cells that composed a cluster of Purkinje cells. In our data set, 54% of the population was composed of burst cells. We computed the population response under the assumption that the membership of a cluster was 54% burst cells. Here, we tested how sensitive the population response was to this membership ratio. The vertical lines indicate saccade onset and offset for all saccades pooled across direction and speed. As the percentage of burst cells in the cluster becomes larger than 70%, or smaller than 50%, the population response no longer returns to baseline at saccade offset.



Extended Data Figure 9 | Gain-field encoding of saccade kinematics in the population response of the Purkinje cells disappeared if the Purkinje cells were organized by their simple-spike activity. **a**, In this analysis we assumed that a collection of 50 Purkinje cells projected onto a single cFN neuron, with the property that all the Purkinje cells shared a similar simple-spike preferred direction. Therefore, the cluster was organized based on the simple-spike properties of the Purkinje cells, not their complex-spike properties. **b**, The population response for saccades made in the direction for which each Purkinje cell showed the largest mean firing rate (simple spikes),

for various saccade peak speeds. The peak population response was not modulated with saccade speed. Error bars are boot-strap-estimated s.e.m. **c**, The population response for saccades made in the direction of maximal modulation. For burst cells, this was the direction for which the Purkinje cell showed the largest mean firing rate, whereas for pause cells, this was the direction associated with the minimum activity (largest pause). The peak population response was not modulated with saccade speed when clusters were organized based on the direction of maximal simple-spike modulation.

η -Secretase processing of APP inhibits neuronal activity in the hippocampus

Michael Willem¹, Sabina Tahirovic², Marc Aurel Busche^{3,4,5}, Saak V. Ovsepian², Magda Chafai⁶, Scherazad Kootar⁶, Daniel Hornburg⁷, Lewis D. B. Evans⁸, Steven Moore⁸, Anna Daria¹, Heike Hampel¹, Veronika Müller¹, Camilla Giudici¹, Brigitte Nuscher¹, Andrea Wenninger-Weinzierl², Elisabeth Kremmer^{2,5,9}, Michael T. Heneka^{10,11}, Dietmar R. Thal¹², Vilmantas Giedraitis¹³, Lars Lannfelt¹³, Ulrike Müller¹⁴, Frederick J. Livesey⁸, Felix Meissner⁷, Jochen Herms², Arthur Konnerth^{4,5}, Hélène Marie⁶ & Christian Haass^{1,2,5}

Alzheimer disease (AD) is characterized by the accumulation of amyloid plaques, which are predominantly composed of amyloid- β peptide¹. Two principal physiological pathways either prevent or promote amyloid- β generation from its precursor, β -amyloid precursor protein (APP), in a competitive manner¹. Although APP processing has been studied in great detail, unknown proteolytic events seem to hinder stoichiometric analyses of APP metabolism *in vivo*². Here we describe a new physiological APP processing pathway, which generates proteolytic fragments capable of inhibiting neuronal activity within the hippocampus. We identify higher molecular mass carboxy-terminal fragments (CTFs) of APP, termed CTF- η , in addition to the long-known CTF- α and CTF- β fragments generated by the α - and β -secretases ADAM10 (a disintegrin and metalloproteinase 10) and BACE1 (β -site APP cleaving enzyme 1), respectively. CTF- η generation is mediated in part by membrane-bound matrix metalloproteinases such as MT5-MMP, referred to as η -secretase activity. η -Secretase cleavage occurs primarily at amino acids 504–505 of APP₆₉₅, releasing a truncated ectodomain. After shedding of this ectodomain, CTF- η is further processed by ADAM10 and BACE1 to release long and short A η peptides (termed A η - α and A η - β). CTFs produced by η -secretase are enriched in dystrophic neurites in an AD mouse model and in human AD brains. Genetic and pharmacological inhibition of BACE1 activity results in robust accumulation of CTF- η and A η - α . In mice treated with a potent BACE1 inhibitor, hippocampal long-term potentiation was reduced. Notably, when recombinant or synthetic A η - α was applied on hippocampal slices *ex vivo*, long-term potentiation was lowered. Furthermore, *in vivo* single-cell two-photon calcium imaging showed that hippocampal neuronal activity was attenuated by A η - α . These findings not only demonstrate a major functionally relevant APP processing pathway, but may also indicate potential translational relevance for therapeutic strategies targeting APP processing.

To identify new proteolytic pathways of APP, we searched for CTFs other than those giving rise to P3 fragment (CTF- α) or amyloid- β (CTF- β)^{3–5}. A new CTF with an approximate molecular mass of 30 kilodaltons (kDa) was revealed, which was recognized by an antibody to the C terminus of APP (Y188) and was absent in the brains of APP-knockout mice⁶ (Fig. 1a and Supplementary Table 1). The molecular mass of the novel CTF suggests an additional physiological cleavage of APP amino-terminal to the known cleavage sites of β - and α -secretases, which we named accordingly η -cleavage of APP (Extended Data

Fig. 1). In the soluble fraction, we detected the N-terminal cleavage product (sAPP- η ; Extended Data Fig. 2), with a molecular mass of approximately 80 kDa that distinguishes it from alternative N-terminal APP fragments described previously^{7–9}. In addition, we observed lower molecular mass soluble peptides (A η), which presumably derived from processing of CTF- η by BACE1 (A η - β) or ADAM10 (A η - α), or alternatively from sAPP- α / β cleavage (Fig. 1b). A η was identified in the soluble fraction of mouse brains as several closely spaced peptides by antibody M3.2 (Fig. 1b), demonstrating that some of these fragments contain the N-terminal part of the amyloid- β domain and probably end at the α -secretase cleavage site. A η fragments were further validated by antibody 9478D directed against an epitope N-terminal to the amyloid- β domain (Fig. 1b). Consistent with η -secretase cleavage of APP in wild-type mice, we observed increased CTF- η and A η production in brain homogenates of APPPS1-21 transgenic mice¹⁰ (Fig. 1c, d). Furthermore, antibody 192swe selectively identified the A η - β species ending at the BACE1 cleavage site (Fig. 1d). Consistent with increased BACE1 cleavage of Swedish mutant APP, only minor amounts of A η - α were detected in this mouse model with antibody 2D8 (Fig. 1d). Physiological η -secretase processing was further confirmed in cerebrospinal fluid (CSF) from humans with and without the Swedish mutation (APP_{swe}) (Fig. 1e). Fivefold higher A η than amyloid- β levels were observed in human CSF (5.33 ± 1.39 times (mean \pm s.d.) higher A η than amyloid- β estimated by 2D8 blot signals; $n = 7$, $P > 0.0001$, Student's *t*-test). Using antibody 192swe, A η - β was selectively detected in the CSF of patients with the Swedish mutation, whereas antibodies 2E9 and 2D8 detected A η - α in all analysed samples (Fig. 1e). Moreover, while these peptides are generated by η -secretase cleavage N-terminal to the amyloid- β domain, they do not reach the γ -secretase site (see mass spectrometric analysis in Fig. 1f), demonstrating that they are different to the previously described N-terminally extended amyloid- β variants^{11,12}. Membrane-bound matrix metalloproteinases such as MT1-MMP and MT5-MMP (also known as MMP14 and MMP24, respectively) were shown to cleave APP *in vitro* at a site consistent with the molecular mass of η -secretase processing products^{13,14}. We therefore produced a neo-epitope-specific antibody (10A8; Extended Data Fig. 2) to identify the η -secretase cleavage site. Antibody 10A8 detected a protein corresponding to sAPP- η with an approximate molecular mass of 80 kDa in mouse brain lysates, which was absent in the APP-knockout mouse brain (Extended Data Fig. 2). Thus, η -secretase cleavage of APP may occur *in vivo* at least in part at amino

¹Biomedical Center (BMC), Ludwig-Maximilians-University Munich, 81377 Munich, Germany. ²German Center for Neurodegenerative Diseases (DZNE) Munich, 81377 Munich, Germany. ³Department of Psychiatry and Psychotherapy, Technische Universität München, 81675 Munich, Germany. ⁴Institute of Neuroscience, Technische Universität München, 80802 Munich, Germany. ⁵Munich Cluster for Systems Neurology (SyNergy), Ludwig-Maximilians-University Munich, 81377 Munich, Germany. ⁶Institut de Pharmacologie Moléculaire et Cellulaire (IPMC), Centre National de la Recherche Scientifique (CNRS), Université de Nice Sophia Antipolis, UMR 7275, 06560 Valbonne, France. ⁷Max Planck Institute of Biochemistry, Martinsried 82152, Germany. ⁸Gurdon Institute, Cambridge Stem Cell Institute & Department of Biochemistry, University of Cambridge, Cambridge CB2 1QN, UK. ⁹Institute of Molecular Immunology, German Research Center for Environmental Health, 81377 Munich, Germany. ¹⁰Department of Neurology, Clinical Neuroscience Unit, University of Bonn, 53127 Bonn, Germany. ¹¹German Center for Neurodegenerative Diseases (DZNE) Bonn, 53175 Bonn, Germany. ¹²Institute of Pathology - Laboratory for Neuropathology, University of Ulm, 89081 Ulm, Germany. ¹³Department of Public Health/Geriatrics, Uppsala University, 751 85 Uppsala, Sweden. ¹⁴Institute for Pharmacy and Molecular Biotechnology IPMB, Functional Genomics, University of Heidelberg, 69120 Heidelberg, Germany.

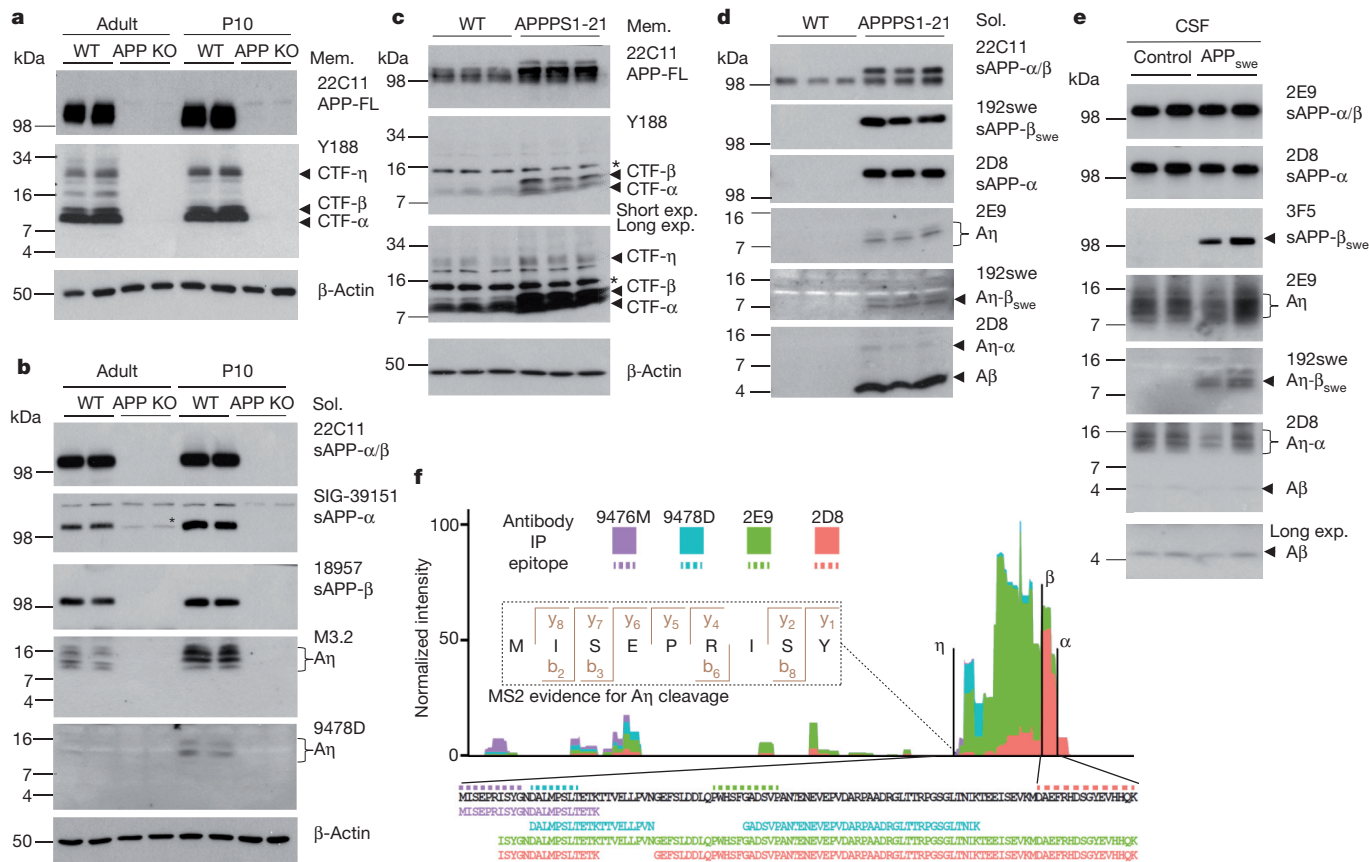


Figure 1 | A novel APP proteolytic processing pathway. **a**, A 30-kDa APP-CTF- η is detected in the brains of adult and postnatal mice (P10, postnatal day 10). APP-FL, full-length APP; KO, knockout; WT, wild type. **b**, An η was detected in the soluble (sol.) fraction of adult and P10 mice by antibody M3.2 and 9478D (antibody 9478D may not be sensitive enough to detect the lower An levels in adult brain). **c**, Higher levels of CTF- η are observed in APPPS1-21 mouse brains as compared to wild-type in the membrane (mem.) fraction. Short and long exposures indicated. Background bands are indicated by

acids 504–505 of APP₆₉₅. To determine the N- and C-terminal cleavage sites of A η peptides, we performed immunoprecipitation with antibodies 9476M, 9478D, 2D8 and 2E9 (Extended Data Fig. 3a, b). Isolated A η peptides were digested with three different proteases to produce several overlapping peptides, and analysed by mass spectrometry. We identified several peptides covering the entire sequence between the cleavage site N504–M505 of APP₆₉₅ starting with the sequence MISEPRISY after the η -secretase cleavage site (Fig. 1f). Mass spectrometry also supports the C-terminal cleavages at the β - and α -secretase sites. After immunoprecipitation with 2D8, fragments of the amyloid- β domain were also observed in smaller amounts alongside the novel A η peptides (Fig. 1f and Extended Data Fig. 3c). As the η -secretase cleavage site at amino acid 505 of APP is consistent with the previously described *in vitro* cleavage sites of APP₆₉₅ by MT1-MMP and MT5-MMP (refs 13, 14), we investigated brains from MT5-MMP- and MT1-MMP-knockout mice^{15,16} (Extended Data Fig. 4). Whereas MT1-MMP knockout had no marked effect on A η - α levels (Extended Data Fig. 4b), the generation of A η - α was reduced in brains from MT5-MMP-knockout mice (Extended Data Fig. 4c). Furthermore, after MT5-MMP overexpression in murine N2a cells, a selective increase in A η - α peptide of approximately 16 kDa was observed (data not shown). Thus, MT5-MMP displays η -secretase activity in intact mouse brains, although the contribution of other η -secretases must be considered.

While investigating protease inhibitors capable of blocking η -secretase, we observed that pharmacological BACE1 inhibition led

asterisks. **d**, Soluble extracts of APPPS1-21 mouse brains contained A η species, as detected by 2E9 antibody. A η - β_{swe} was selectively detected by antibody 192swe. **e**, A η and amyloid- β (A β) were readily detectable in human CSF by antibody 2D8. Antibody 2E9 allowed the detection of A η in all samples, whereas 192swe specifically detected A η - β_{swe} . **f**, Mass spectrometric analysis of A η . Peptide intensities were summed per amino acid residue and plotted in relation to each other.

to a pronounced accumulation of the long A η - α species in Chinese hamster ovary (CHO) 7PA2 cells (Fig. 2a). This indicates that after blockade of β -secretase activity, processing by α -secretase leads to enhanced production of the long A η - α species at the expense of the shorter BACE1-generated A η - β . Similarly, BACE1 inhibition also led to an accumulation of endogenous CTF- η and enhanced production of endogenous A η - α in primary mouse hippocampal neurons (Fig. 2b, c), as well as human neurons differentiated from embryonic pluripotent stem cells (H9 cells¹⁷; Fig. 2d–g). Furthermore, in human neurons we not only detected a 65% increase (Fig. 2g) in the slightly longer A η - α species after BACE inhibition, but also a concomitant decrease in A η - β peptides (Fig. 2e, top). Importantly, η -secretase processing significantly exceeded amyloidogenic processing (9.5 ± 1.87 times A η compared to amyloid- β estimated for human neurons; $n = 8$, $P > 0.001$, Student's t -test). Pharmacological intervention with a BACE1 inhibitor *in vivo* caused a clear and time-dependent increase in CTF- η and A η - α levels in APP_{V717I} transgenic mice¹⁸ (Fig. 2h), which was fully reversible within 24 h after administration. In agreement, increased CTF- η and A η - α levels were also observed in a BACE1-knockout mouse¹⁹ *in vivo* (Fig. 2i).

To identify a potential contribution of η -secretase processing to AD pathology, immunohistochemical analyses were performed with brains derived from six-month-old APPPS1-21 (ref. 20) mice (Extended Data Fig. 5a). This revealed co-labelling of antibody Y188 with antibody 2E9 in dystrophic neurites. No signal for A η peptides was obtained in plaque cores in which aggregated amyloid- β was

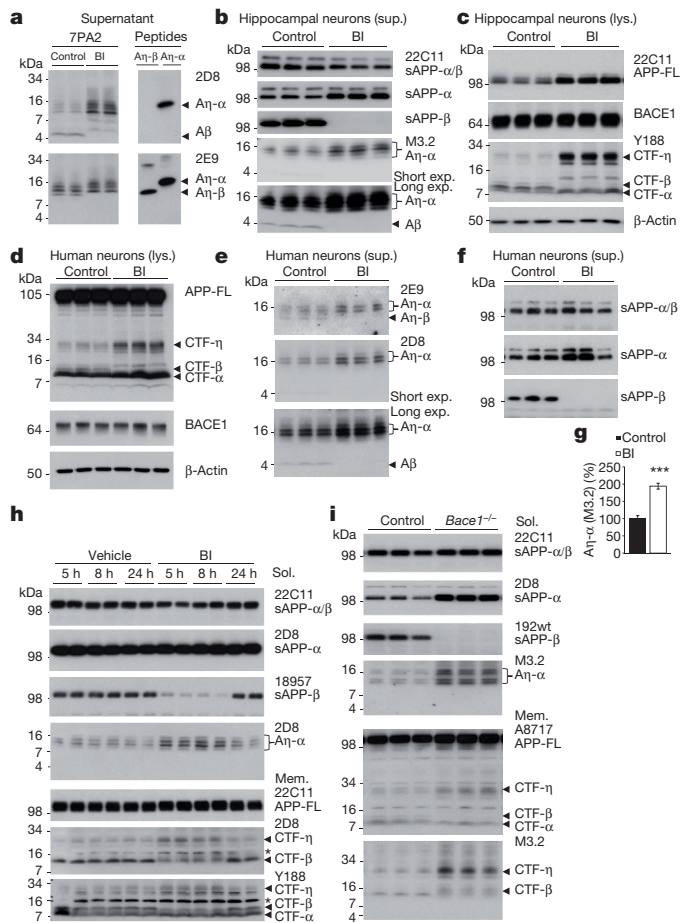


Figure 2 | Inhibition of BACE1 results in increased levels of CTF- η and A η - α . **a**, Conditioned media from 7PA2 cells treated with or without a BACE1 inhibitor (BI) was co-migrated with synthetic A η - β and A η - α peptides and immunoblotted with 2D8 and 2E9 antibodies. **b–f**, After BACE inhibition in mouse hippocampal neurons and human neurons, a reduction of sAPP- β was accompanied by a strong increase in endogenous A η - α and CTF- η levels. Lys., lysate; sup., supernatant. **g**, Quantification of intensities for 2D8 signals in **e** ($n = 8$; *** $P < 0.001$, Student's t -test). **h**, BACE1 inhibition *in vivo* resulted in enhanced production of A η - α species in APP^{V717I} mice. Background bands are indicated by asterisks. **i**, Western blot analysis of soluble extracts of P10 *Bace1*^{-/-} mouse brains revealed a marked increase in A η - α peptides as compared to controls.

detected by 6E10 staining (Extended Data Fig. 5a). Similar data were obtained in human AD brains (Extended Data Fig. 6). To verify accumulation of CTF- η in dystrophic neurites, we used laser capture microdissection (LCM). Western blot analysis revealed not only CTF- β and CTF- α , but also CTF- η within the halo, but not within the plaque core area or regions devoid of plaques (Extended Data Fig. 5b). As expected, amyloid- β was observed within the plaque core as well as in the surrounding halo (Extended Data Fig. 5b).

Since the cleavage products of η -secretase APP processing accumulate after BACE1 inhibition and are enriched in dystrophic neurites, we examined whether soluble A η peptides interfere with neuronal function, similar to soluble amyloid- β oligomers¹. Long-term potentiation (LTP) is considered as a synaptic correlate of memory, and is widely used as a model for investigating the neurotoxic effects of amyloid- β oligomers on synaptic function^{21–23}. A single oral dose of the BACE1 inhibitor SCH1682496 increased the CTF- η levels, and almost doubled the A η - α level in soluble brain extracts prepared from animals 3 h after treatment (Extended Data Fig. 7). This was accompanied by a significant reduction of hippocampal LTP (Fig. 3a, b). These findings may suggest an involvement of A η - α in LTP deficit

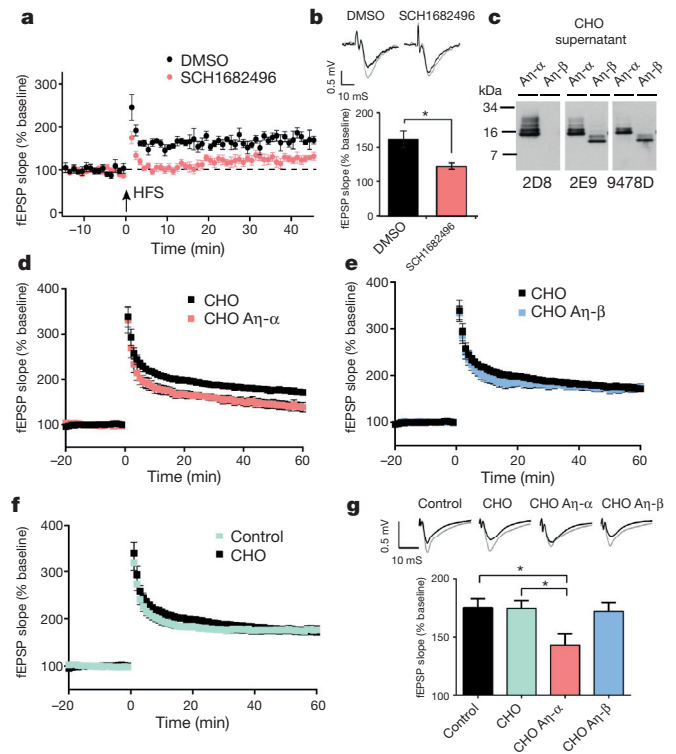


Figure 3 | A η - α impairs hippocampal LTP. **a**, Pharmacological BACE1 inhibition with SCH1682496 lowers hippocampal LTP. DMSO, dimethylsulfoxide; fEPSP, field excitatory postsynaptic potential; HFS, high-frequency stimulation. **b**, Representative fEPSPs recorded in CA1 area before and 45 min after tetanization of Schaffer collaterals (top), with summary plot of the effects of the inhibitor (SCH1682496) and vehicle (DMSO) on fEPSP slopes in all examined groups in **a** ($n = 9$). **c**, Soluble A η - α and A η - β peptides were expressed in CHO cells and blotted with 2D8, 2E9 and 9478D antibodies. **d–f**, A η - α ($n = 9$) (**d**), but not A η - β ($n = 7$) (**e**), conditioned media from untransfected cells (CHO; $n = 13$) or artificial cerebrospinal fluid (ACSF) (control; $n = 15$) (**f**) significantly inhibited LTP. **g**, Summary graph of LTP magnitudes calculated 45–60 min after high-frequency stimulation from graphs in **d–f** with statistical analysis (* $P < 0.05$; one-way analysis of variance (ANOVA) and post hoc Bonferroni test); error bars represent s.e.m. For each condition, sample fEPSP traces pre-LTP (black) and 45–60 min post-LTP (grey) induction are shown (top).

under acute blockade of β -secretase activity. To validate the potential effects of A η peptides on synaptic transmission and plasticity directly, we expressed A η - β and A η - α in CHO cells (Fig. 3c). Concentrated conditioned media was further enriched for A η by size-exclusion chromatography (SEC) and applied to hippocampal slices before LTP induction in CA1 pyramidal neurons. Neither A η - β nor A η - α influenced the baseline synaptic transmission (Extended Data Fig. 8). Comparison of LTP 60 min after its induction in the presence of A η - β or A η - α with control conditions (Fig. 3d–g) revealed that A η - α lowered the LTP to a degree comparable to synthetic amyloid- β dimers²³ (Extended Data Fig. 9a, b), while truncated A η - β had no effect (Fig. 3e, g). In support of this observation, synthetic A η - α reduced LTP to a similar extent at a concentration of 100 nM (Extended Data Fig. 9c, d). To examine the direct effects of A η peptides on neuronal activity *in vivo*, we used two-photon calcium imaging at single-cell resolution^{24,25}. Figure 4a–d illustrates results from experiments in which the activity of neurons in the CA1 pyramidal cell layer was monitored before and after superfusion of the exposed hippocampus with A η peptides or the respective control peptides. A η - α strongly suppressed the activity of hippocampal neurons *in vivo*, an effect not observed with A η - β or the control peptide (Fig. 4a–d and Extended Data Fig. 10). By using local application of synthetic A η - α to hippocampal neurons, we demonstrate that the inhibitory effect of A η - α on neurons was readily reversible after

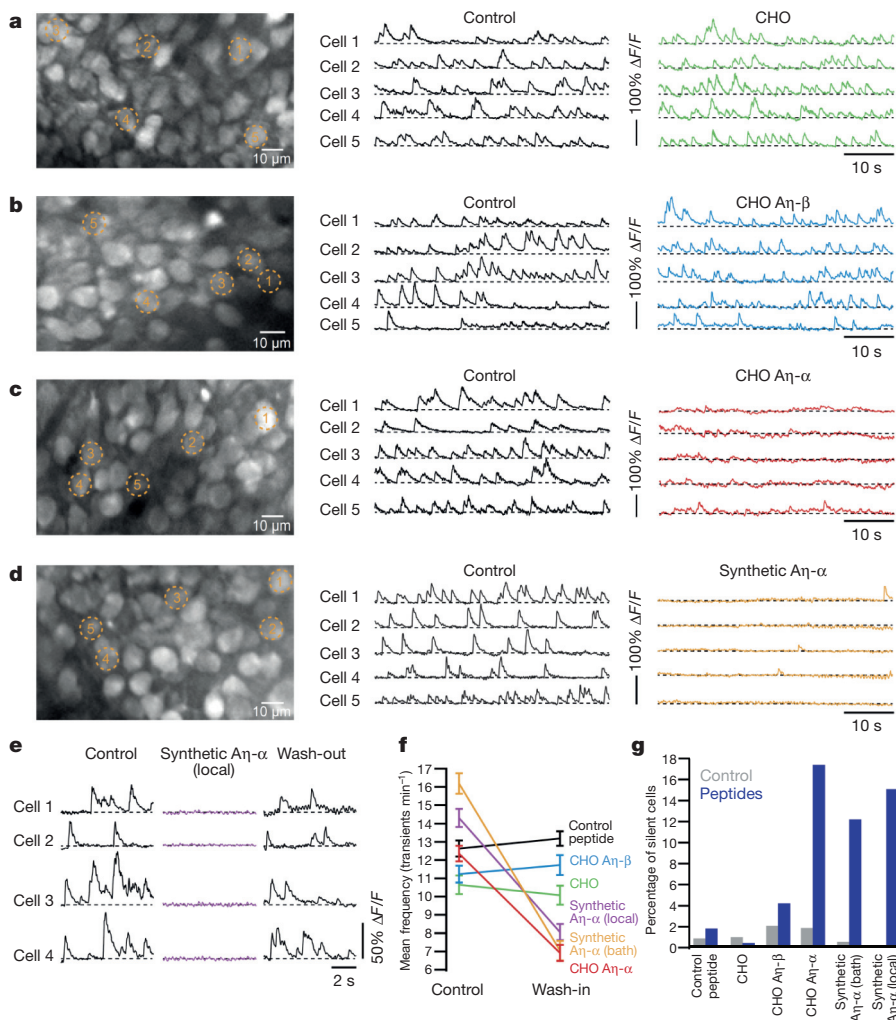


Figure 4 | **Aη-α reduces neuronal activity *in vivo*.** **a–d**, Left, *in vivo* images of CA1 hippocampal neurons labelled with the fluorescent calcium indicator fluo-8AM. Middle and right, calcium transients of five representative neurons, marked in the corresponding left panels, before and during bath-application of Aη or CHO conditioned media. **e**, Calcium transients in hippocampal neurons before, during and after local application of synthetic Aη-α. **f, g**, Summary results of the changes in the average rates of calcium transients (**f**), and in the fractions of silent neurons (**g**); for statistics see Supplementary Tables 2 and 3. Error bars represent s.e.m.

washout (Fig. 4e). A summary of the results from all experiments is shown in Fig. 4f, g.

We have identified a new APP processing pathway that exceeds amyloidogenic processing. Similar to amyloid-β production²⁶, the alternative proteolytic processing pathway occurs under physiological conditions but may be altered during AD pathogenesis. Accumulation of η-secretase²⁷ and CTF-η within dystrophic neurites in close vicinity to neuritic plaques may also support its potential contribution to AD pathology. However, all APP and presenilin-associated familial AD mutations affect amyloid-β production and aggregation (reviewed in ref. 26), whereas the Icelandic mutation APP_{A673T} prevents AD and dementia by moderately reducing amyloid-β production²⁸. Indeed, the Swedish mutation decreased Aη-α by strongly enhancing BACE1-mediated APP processing. However, η-secretase accumulation²⁷ and its activity near amyloid plaques may suggest that η-secretase stimulation by amyloid-β could be a downstream effector within the amyloid cascade. Although Aη may be involved in the modulation of neuronal activity and synaptic plasticity, the differential bioactivity of recombinant Aη-α and Aη-β is currently unclear. One may speculate that the longer Aη-α peptide is more stable owing to unknown post-translational modifications. This would be consistent with our observation that, in contrast to cell-produced Aη-β, 100 nM of synthetic Aη-β inhibits LTP (data not shown). Finally, it is important to note that low-*n* amyloid-β oligomer preparations from 7PA2 supernatants contain considerable amounts of Aη (data not shown). Thus, the previously observed inhibition of LTP with such fractions may also be attributed to the presence of Aη. Our findings may also be considered in the context of continuing clinical trials with BACE1 inhibitors. Together with the identification of numerous

brain-specific BACE1 substrates^{29,30}, our data indicate that therapeutic inhibition of BACE1 activity requires careful titration to prevent unwanted adverse effects at several levels.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 27 March; accepted 25 June 2015.

Published online 31 August 2015.

- Haass, C. & Selkoe, D. J. Soluble protein oligomers in neurodegeneration: lessons from the Alzheimer's amyloid β-peptide. *Nature Rev. Mol. Cell Biol.* **8**, 101–112 (2007).
- Dobrowolska, J. A. *et al.* CNS amyloid-β, soluble APP-α and -β kinetics during BACE inhibition. *J. Neurosci.* **34**, 8336–8346 (2014).
- Haass, C., Koo, E. H., Mellon, A., Hung, A. Y. & Selkoe, D. J. Targeting of cell-surface β-amyloid precursor protein to lysosomes: alternative processing into amyloid-bearing fragments. *Nature* **357**, 500–503 (1992).
- Estus, S. *et al.* Potentially amyloidogenic, carboxyl-terminal derivatives of the amyloid protein precursor. *Science* **255**, 726–728 (1992).
- Haass, C. *et al.* β-Amyloid peptide and a 3-kDa fragment are derived by distinct cellular mechanisms. *J. Biol. Chem.* **268**, 3021–3024 (1993).
- Müller, U. *et al.* Behavioral and anatomical deficits in mice homozygous for a modified β-amyloid precursor protein gene. *Cell* **79**, 755–765 (1994).
- Nikolaev, A., McLaughlin, T., O'Leary, D. D. & Tessier-Lavigne, M. APP binds DR6 to trigger axon pruning and neuron death via distinct caspases. *Nature* **457**, 981–989 (2009).
- Jefferson, T. *et al.* Metalloprotease meprin β generates nontoxic N-terminal amyloid precursor protein fragments *in vivo*. *J. Biol. Chem.* **286**, 27741–27750 (2011).
- Vella, L. J. & Cappai, R. Identification of a novel amyloid precursor protein processing pathway that generates secreted N-terminal fragments. *FASEB J.* **26**, 2930–2940 (2012).
- Radde, R. *et al.* β42-driven cerebral amyloidosis in transgenic mice reveals early and robust pathology. *EMBO Rep.* **7**, 940–946 (2006).

11. Portelius, E. *et al.* Mass spectrometric characterization of amyloid- β species in the 7PA2 cell model of Alzheimer's disease. *J. Alzheimers Dis.* **33**, 85–93 (2013).
12. Welzel, A. T. *et al.* Secreted amyloid β -proteins in a cell culture model include N-terminally extended peptides that impair synaptic plasticity. *Biochemistry* **53**, 3908–3921 (2014).
13. Higashi, S. & Miyazaki, K. Novel processing of β -amyloid precursor protein catalyzed by membrane type 1 matrix metalloproteinase releases a fragment lacking the inhibitor domain against gelatinase A. *Biochemistry* **42**, 6514–6526 (2003).
14. Ahmad, M. *et al.* Cleavage of amyloid- β precursor protein (APP) by membrane-type matrix metalloproteinases. *J. Biochem.* **139**, 517–526 (2006).
15. Folgueras, A. R. *et al.* Metalloproteinase MT5-MMP is an essential modulator of neuro-immune interactions in thermal pain stimulation. *Proc. Natl Acad. Sci. USA* **106**, 16451–16456 (2009).
16. Zhou, Z. *et al.* Impaired endochondral ossification and angiogenesis in mice deficient in membrane-type matrix metalloproteinase 1. *Proc. Natl Acad. Sci. USA* **97**, 4052–4057 (2000).
17. Shi, Y., Kirwan, P. & Livesey, F. J. Directed differentiation of human pluripotent stem cells to cerebral cortex neurons and neural networks. *Nature Protocols* **7**, 1836–1846 (2012).
18. Moechars, D. *et al.* Early phenotypic changes in transgenic mice that overexpress different mutants of amyloid precursor protein in brain. *J. Biol. Chem.* **274**, 6483–6492 (1999).
19. Cai, H. *et al.* BACE1 is the major β -secretase for generation of A β peptides by neurons. *Nature Neurosci.* **4**, 233–234 (2001).
20. Radde, R. *et al.* A β 42-driven cerebral amyloidosis in transgenic mice reveals early and robust pathology. *EMBO Rep.* **7**, 940–946 (2006).
21. Walsh, D. M. *et al.* Naturally secreted oligomers of amyloid β protein potently inhibit hippocampal long-term potentiation *in vivo*. *Nature* **416**, 535–539 (2002).
22. Shankar, G. M. *et al.* Natural oligomers of the Alzheimer amyloid- β protein induce reversible synapse loss by modulating an NMDA-type glutamate receptor-dependent signaling pathway. *J. Neurosci.* **27**, 2866–2875 (2007).
23. Shankar, G. M. *et al.* Amyloid- β protein dimers isolated directly from Alzheimer's brains impair synaptic plasticity and memory. *Nature Med.* **14**, 837–842 (2008).
24. Busche, M. A. *et al.* Clusters of hyperactive neurons near amyloid plaques in a mouse model of Alzheimer's disease. *Science* **321**, 1686–1689 (2008).
25. Busche, M. A. *et al.* Critical role of soluble amyloid- β for early hippocampal hyperactivity in a mouse model of Alzheimer's disease. *Proc. Natl Acad. Sci. USA* **109**, 8740–8745 (2012).
26. Haass, C. Take five—BACE and the γ -secretase quartet conduct Alzheimer's amyloid β -peptide generation. *EMBO J.* **23**, 483–488 (2004).
27. Sekine-Aizawa, Y. *et al.* Matrix metalloproteinase (MMP) system in brain: identification and characterization of brain-specific MMP highly expressed in cerebellum. *Eur. J. Neurosci.* **13**, 935–948 (2001).
28. Jonsson, T. *et al.* A mutation in APP protects against Alzheimer's disease and age-related cognitive decline. *Nature* **488**, 96–99 (2012).
29. Kuhn, P. H. *et al.* Secretome protein enrichment identifies physiological BACE1 protease substrates in neurons. *EMBO J.* **31**, 3157–3168 (2012).
30. Zhou, L. *et al.* The neural cell adhesion molecules L1 and CHL1 are cleaved by BACE1 protease *in vivo*. *J. Biol. Chem.* **287**, 25927–25940 (2012).

Supplementary Information is available in the online version of the paper.

Acknowledgements The authors thank S. Lammich, N. Exner and H. Steiner for critical comments. We thank A. Sülzen, N. Astola, S. Diederich, E. Grieflinger and J. Gobbert for technical help. The APPPS1-21 colony was established from a breeding pair provided by M. Jucker. *MT1-MMP*^{−/−} mouse brains were obtained from Z. Zhou. *MT5-MMP*^{−/−} mouse brains were obtained from I. Farinas. We thank H. Jacobsen for the BACE1 inhibitor RO5508887. This work was supported by the European Research Council under the European Union's Seventh Framework Program (FP7/2007–2013)/ERC grant agreement no. 321366-Amyloid (advanced grant to C.H.). The work of D.R.T. was supported by AFI (grant 13803). The research leading to these results has received funding (F.M. and D.H.) from the European Research Council under the European Union's Seventh Framework Programme (FP7/2007–2013)/ERC grant agreement no. 318987 [TOPAG]. We thank J. Cox and M. Mann for critical discussions and the mass spectrometry infrastructure. We also acknowledge support by grants from Deutsche Forschungsgemeinschaft (MU 1457/9-1, 9-2 to U.M.) and the ERA-Net Neuron (01EW1305A to U.M.). Further support came from the ATIP/AVENIR program (Centre national de la recherche scientifique, CNRS) to H.M.; the French Fondation pour la Coopération Scientifique – Plan Alzheimer (Senior Innovative Grant 2010) to M.C. and H.M., and the French Government (National Research Agency, ANR) through the “Investments for the Future” LABEX SIGNALIFE: program reference ANR-11-LABX-0028-01 to S.K. M.A.B. was supported by the Langmatz Stiftung. F.J.L. is a Wellcome Trust Investigator. *In vivo* BACE1 inhibition experiments with APP transgenic mice were performed together with reMYND (Bio-Incubator, 3001 Leuven-Heverlee, Belgium).

Author Contributions M.W. and C.H. designed the study and interpreted the results. M.W. generated all biochemical data together with H.H., V.M., B.N. and C.G. S.T., supported by A.W.-W., provided primary neuronal cultures, performed and analysed immunohistological stainings, and together with A.D. performed LCM. D.R.T. provided and analysed human brain sections. M.T.H. provided CSF samples. U.M. provided APP-knockout mice. E.K. produced new monoclonal antibodies. D.H. and F.M. designed and conducted mass spectrometry and data analysis. L.D.B.E., S.M. and F.J.L. carried out BACE1 inhibition of human neurons. H.M. together with M.C. and S.K. performed all electrophysiological recordings (LTP) *in vitro* and analysis in relation to application of peptides. S.V.O. and J.H. performed all electrophysiological recordings (LTP) *in vitro* in relation to the BACE1 inhibitor tests. M.A.B. and A.K. performed all Ca²⁺-imaging experiments *in vivo* and analysis. M.W. and C.H. wrote the manuscript with input from the other authors. Correspondence and requests for materials should be addressed to M.W. and C.H.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare competing financial interests: details are available in the online version of the paper. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to M.W. (michael.willem@mail03.med.uni-muenchen.de) or C.H. (christian.haass@mail03.med.uni-muenchen.de).

METHODS

Cell culture. Mycoplasma-free CHO and 7PA2 (ref. 31) cells (gift from E. Koo) were grown in DMEM/F12 (Thermo Scientific) supplemented with 10% FCS (Thermo Scientific) plus penicillin/streptomycin and non-essential amino acids (PAA) in a humid incubator with 5% CO₂ at a temperature of 37 °C. For inhibitor treatment, cell culture medium was replaced with fresh, pre-warmed serum free medium (OPTIMEM; Thermo Scientific) supplemented with inhibitors or DMSO as vehicle control. Treatment was initiated when cells reached 90–100% confluency and conditioned media were collected after 20–24 h. Supernatants were cleared by centrifugation (10 min, 5,500g at 4 °C). To obtain cell lysates, cell monolayers were washed once with ice-cold PBS and detached in 1 ml PBS using a cell scraper. The cell suspension was pelleted by centrifugation (5 min, 1,000g at 4 °C) and lysed with RIPA buffer (20 mM sodium citrate, pH 6.4, 1 mM EDTA, 1% Triton X-100 in ddH₂O) supplemented with Protease Inhibitor Cocktail (Sigma-Aldrich). The protein concentration of lysates was determined using the Uptima BC Assay Protein Quantitation kit (Interchim).

Primary cell culture. Hippocampal neurons were isolated from embryonic day 18 CD rats (Charles River) as described previously³². Dissociated neurons were plated at 17,700 cells cm⁻² onto 6-cm dishes coated with poly-L-lysine (1 mg ml⁻¹; Sigma) and cultured in Neurobasal medium supplemented with 2% B27 and 0.5 mM L-glutamine (all from Invitrogen). Hippocampal cultures were maintained in a humidified 5% CO₂ incubator at 37 °C. For inhibitor treatment, DIV16 culture medium was replaced with fresh, pre-equilibrated N2 medium (supplemented with 20% of 4 days conditioned N2 medium from pure primary cultured astrocytes) to which inhibitors or DMSO as vehicle control were added.

Generation and BACE1 inhibition of cerebral cortex neurons induced from human embryonic stem cells. Cell lines in this study were H9 ES (WiCell Research Institute)³³. Pluripotent cells were cultured on mouse embryonic fibroblasts (GlobalStem) in DMEM/F12 containing 20% (v/v) KSR, 100 µM non-essential amino acids, 100 µM 2-mercaptoethanol, 50 U ml⁻¹ penicillin and 50 mg ml⁻¹ streptomycin (Life Technologies) and 10 ng ml⁻¹ FGF2. Directed differentiation of human embryonic stem cells to cerebral cortex neurons was carried out as described^{17,34}. Human neurons (75 days after induction) were treated with 1 µM β -secretase inhibitor LY2886721 (Selleck) dissolved in DMSO (20 mM stock). Vehicle-only control assays were performed using DMSO. The compound was applied twice at 48-h intervals. Extracellular media was collected before drug addition and at subsequent 48-h intervals. Neurons were collected after 4 days of treatment using 0.5 mM EDTA in PBS.

Transgenic mice, animal care and animal handling. *Bace1*^{-/-} and APPS1-21 mice were described before^{10,19} and were bred for this study in a Bl6C57/J background. All treatments were approved by the local committee for animal use and were performed in accordance to state and federal regulations (license number KVR-I/221-TA116/09). Mice had access to pre-filtered sterile water and standard mouse chow (Ssniff Ms-H, Ssniff Spezialdiäten GmbH, Soest, Germany) *ad libitum* and were housed under a reversed day–night rhythm in IVC System Type II L-cages (528 cm²) equipped with solid floors and a layer of bedding, in accordance to local legislation on animal welfare.

BACE1 inhibitor treatment. Randomized APP_{V7171} (ref. 18) mice were treated with vehicle or with the inhibitor RO5508887 provided by Hoffmann-La Roche³⁵. The groups of treated mice were blinded to the examiner and uncoded at the end of the experiments. Three-month-old heterozygous female transgenic mice in mixed FVB/N × C57Bl/6J background expressing human APP_{V7171} (ref. 18) were used for BACE1 inhibition studies. Gavage mediated administration of BACE1 inhibitor (90 mg kg⁻¹, 14.06 ml kg⁻¹) or vehicle (14.06 ml kg⁻¹) was performed once³⁵. The BACE1 inhibitor was diluted in 5% ethanol (Merck) and 10% solutol (Sigma-Aldrich) in sterile water (Baxter). Animals were sacrificed after 5, 8 and 24 h. Mice were anaesthetized with 3.5 µl per gram body weight of a mixture of ketamine (115 mg ml⁻¹ ketamine hydrochloride, Eurovet), xylazine 2% (23.32 mg ml⁻¹ xylazine hydrochloride, VMD Arendonk), atropine (0.50 mg ml⁻¹ atropine sulphate, Sterop) and saline (8.5:2.5, v/v/v/v). For brain preparation, mice were flushed *trans*-cardially with ice-cold saline (3.5 ml min⁻¹, 3 min). The brain was removed from the cranium and dissected into left and right hemiforebrain, brainstem, cerebellum and olfactory bulb. The brain structures were promptly immersed in liquid nitrogen and stored at -80 °C. Different tissues (kidneys, spleen, liver, stomach, gut, lungs and heart) were examined and checked for gross abnormalities. No obvious abnormalities were observed in any of the treatment groups.

Preparation of protein extracts from brain. Brains were removed from the cranium and dissected into left and right hemispheres. Brain tissue was snap-frozen in liquid nitrogen and stored at -80 °C. Soluble proteins were extracted with DEA buffer (50 mM NaCl, 0.2% diethylamine, pH 10, plus protease inhibitor (P8340, Sigma-Aldrich))³⁶, membrane proteins were extracted with RIPA buffer (20 mM Tris-HCl, pH 7.5, 150 mM NaCl, 1 mM EDTA, 1 mM EGTA, 1% NP-40, 1% sodium deoxycholate, 2.5 mM sodium pyrophosphate

plus protease inhibitor) or applying a membrane preparation protocol as described before³⁷.

Protein analysis. Proteins were separated under denaturing conditions using discontinuous SDS-PAGE. Equal amounts of proteins denatured in Laemmli buffer were loaded onto the gel and 10 µl of the SeeBlue Plus2 Prestained Standard (Invitrogen) served as molecular mass marker. Electrophoresis was performed in Tris-glycine buffer (25 mM Tris, 190 mM glycine in ddH₂O) using the Mini-PROTEAN system (BIORAD) on activated PVDF membranes. Low molecular mass proteins (<16 kDa) were separated using precast gradient Tricine Protein Gels (10–20%, 1 mm, Novex) in Tris-tricine buffer using the XCell SureLock Mini-Cell system (Novex). After separation by SDS-PAGE, proteins were transferred onto membranes using the tank/wet Mini Trans-Blot cell system (BIORAD). CTFs, A η and amyloid- β were detected after transfer on Nitrocellulose membranes (Protran BA85; GE Healthcare), while other proteins were blotted on PVDF (Immobilon-P, Merck Millipore). As size markers for A η synthetic peptides A η - β (92 amino acids; 1-MISEPRISYGNDALMPSLTETKT TVELLPVNGEFLDLDLPWHSFGADSV PANTENEVEVPDARPAADRGLTT RPSGLTNIKTEEISEVKM-92) and the slightly longer A η - α (108 amino acids; 1-MISEPRISYGNDALMPSLTETKT TVELLPVNGEFLDLDLPWHSFGADSV PANTENEVEVPDARPAADRGLTTRPSGLTNIKTEEISEVKMDAEFRHDSG YEVHHQK-108) were obtained from Peptide Speciality Laboratories. After completion of the transfer and before blocking, proteins transferred to nitrocellulose membranes were additionally denatured by boiling the membrane in PBS (140 mM NaCl, 10 mM Na₂HPO₄, 1.75 mM KH₂PO₄, 2.7 mM KCl in ddH₂O, pH 7.4) for 5 min. After cooling to room temperature, nitrocellulose membranes as well as the PVDF membranes were blocked in I-Block solution (0.2% Tropix I-Block (Applied Biosystems), 0.1% Tween20 in PBS) for 1 h at room temperature or overnight at 4 °C (with agitation). Transferred proteins were detected using immunodetection and enhanced chemiluminescence (ECL). First, blocked membranes were incubated with primary antibodies diluted in I-Block solution overnight at 4 °C (with agitation). After removal of the antibody, membranes were washed three times in TBS-T buffer (10 min each, at room temperature, with agitation; 140 mM NaCl, 2.68 mM KCl, 24.76 mM Tris, 0.3% Triton X-100 in ddH₂O, pH 7.6) and subsequently incubated with a horseradish-peroxidase-coupled secondary antibody (obtained from Promega or Santa Cruz). Secondary antibodies were diluted in I-Block solution and membranes were incubated for 1 h at room temperature (with agitation) followed by three washes in TBS-T. For ECL detection, membranes were incubated with horseradish peroxidase substrate (ECL, GE Healthcare or ECL Plus, Thermo Scientific) for 1 min at room temperature and signals were captured with X-ray films (Super RX Medical X-Ray, Fujifilm), which were subsequently developed using an automated film developer (CAWOMAT 2000 IR, CAWO). Quantitation of protein was conducted using ImageJ software. Ratios were obtained from signals on the same film for A η over amyloid- β . Quantitative data were analysed statistically by using a two-tailed Student's *t*-test.

Molecular cloning and transfection. For the expression of A η - α and A η - β in CHO cells, the complementary DNAs of the respective fragments were amplified by PCR and subcloned into the pSecTag2A (Invitrogen) vector that features an N-terminal secretion signal. CHO cells were cultured in DMEM with 10% FCS and non-essential amino acids. Transfections were carried out using Lipofectamine 2000 (Invitrogen) according to the manufacturer's instructions.

Mass spectrometry analysis of samples. Beads with immunoprecipitated peptides were resuspended in ddH₂O and reduced with 10 mM dithiothreitol followed by alkylation with 55 mM 2-chloroacetamide. Samples were divided into three parts and digested with either 1 µg of trypsin (1.2 M urea, 0.4 M thiourea and 50 mM ammonium bicarbonate), LysC (5 M urea, 1.7 M thiourea and 50 mM ammonium bicarbonate) or chymotrypsin (0.3 M urea, 0.1 M thiourea and 50 mM ammonium bicarbonate). To increase the sequence coverage further, partially cleaved peptides were generated by digesting for 5, 10, 20, 40, 60, 120, 180 and 720 min. Samples from all time points of a respective protease were pooled and desalted on stage tips³⁸.

For liquid chromatography–tandem mass spectrometry (LC–MS/MS), peptides were separated on a Thermo Scientific EASY-nLC 1000 HPLC system (Thermo Fisher Scientific) and in-house packed columns (75 µm inner diameter, 20 cm length, 1.9 µm C18 particles (Dr. Maisch GmbH). The peptide mixture was loaded in buffer A (0.5% formic acid) and separated with a gradient from 10% to 60% buffer B (80% acetonitrile, 0.5% formic acid) within 40 min at 250 nl min⁻¹ at a column temperature of 50 °C. A Quadrupole Orbitrap mass spectrometer³⁹ (Q Exactive, Thermo Fisher Scientific) was coupled to the HPLC system via a nano electrospray source. We used data-dependent acquisition with a survey scan range of 300 to 1,650 *m/z*, at a resolution of 60,000 *m/z* and selected up to five most abundant features with a charge state ≥ 2 for HCD fragmentation⁴⁰ at a normalized collision energy of 27 and a resolution of 15,000 at *m/z* 200. To limit repeated

sequencing, dynamic exclusion of sequenced peptides was set to 20 s. Thresholds for ion injection time and ion target values were set to 20 ms and 3×10^6 for the survey scans, and 120 ms and 1×10^5 for the MS/MS scans. Data were acquired using the Xcalibur software (Thermo Scientific).

Data analysis. To process mass spectrometry raw files, we used the MaxQuant software (v1.5.2.16)⁴¹. We used the Andromeda search engine⁴², which is integrated into MaxQuant, to search MS/MS spectra against the APP₆₉₅ and 247 common contaminating proteins⁴². We set enzyme specificity to unspecific to detect novel cleavage sites and set a peptide search length from 7 to 40 amino acids. A false discovery rate cutoff of 1% was applied at the peptide level. For data visualization we used R⁴³. Identified peptides were mapped to APP₆₉₅. To display quantitative evidence for overlapping peptides, intensities of identified peptides were summed and plotted per amino acid residue. The data of the individual immunoprecipitation and mass spectrometry analyses are depicted in Extended Data Fig. 3.

Human CSF samples. Human CSF samples collected at the Department of Neurology Outpatient unit for neurodegenerative disease (KBfZ) of the University of Bonn were obtained by lumbar puncture at position L3, centrifuged and divided in small aliquots. For further analysis, samples were stored at -80°C . Turbid or blood-contaminated samples were excluded from analysis. Use of these samples for research purposes has been consented by all patients according to the ethical committee requirements of the University of Bonn Ethical committee and approval number 279/10. For the analysis of APP_{swc} carriers with antibodies 192swc (ref. 44) lumbar CSF was obtained from family members. Tubes with CSF were stored at -70°C until analysis. The clinical diagnosis of probable AD was based on NINCDS-ADRDA criteria⁴⁵. The diagnosis of AD was confirmed by neuropathological examination of the brain of one deceased mutation-carrier^{46,47}. This study was approved by research ethics committee at the Uppsala University Hospital (Dnr 048-2005).

Neuropathology and immunohistochemistry. Use of brain samples for research purposes has been consented by all patients according to the ethical committee requirements of the University of Ulm Ethical committee and approval number 54/08. Braak-NFT stages⁴⁸, and CERAD⁴⁹ scores for neuritic plaques were used to determine the degree of AD pathology according to the NIA-AA guidelines⁵⁰. Consecutive paraffin sections from the human medial lobe were stained with 22C11, 9476M and 9478D. Primary antibodies were detected with biotinylated anti-mouse and anti-rabbit IgG secondary antibodies and visualized with avidin-biotin-complex (ABC-Kit, Vector Laboratories) and diaminobenzidine-HCl (DAB). The sections were counterstained with haematoxylin. Positive and negative controls were performed. 9476M and 9478D stainings were assessed in 10 control and 10 AD patient cases.

For double immunofluorescence analysis of APPPS1-21 brain sections, 6-month-old mice were killed by CO₂ inhalation according to animal handling laws. Brains were dissected and fixed with 4% paraformaldehyde in 0.1 M PBS, pH 7.4 for 48 h. For immunohistochemistry, 25- μm -thick sagittal mouse brain cryosections were treated with 10 mM sodium citrate, pH 6 at 95°C for 20 min, washed with 0.5% Triton X-100 in PBS, blocked with 5% goat serum (Invitrogen) and 0.5% Triton X-100 in PBS for 1 h and subsequently incubated overnight with primary antibodies diluted in blocking solution. Primary antibodies were used as listed in Supplementary Table 1. DAPI was used to counterstain nuclei. Signals were visualized using fluorescently labelled secondary antibodies. Confocal images were acquired using a Plan-Apochromat 25 \times /0.8 oil differential interference contrast objective on a LSM 710 confocal microscope (Zeiss) in sequential scanning mode using ZEN 2011 software package (black edition, Zeiss).

LCM of plaque enriched brain material. For laser capture microdissection of plaque cores and halos, 10-, 11-, 14-, 16- and 24-month-old transgenic APPPS1-21 mice were used according to a previously published protocol⁵¹ with slight modifications. Mice brains were dissected and immediately frozen on crushed dry ice. Ten-micrometre-thick sagittal sections were cut using a Microm HM 560 cryostat (Thermo Scientific), mounted on frame slides containing a 1.4 μm polyethylene terephthalate membrane (Leica Microsystems) and subsequently stained or stored at -80°C for later usage. Staining was performed as follows: brain sections were thawed briefly at room temperature, fixed with 75% ethanol for 1 min, stained with 0.05% Thioflavin-S for 5 min, washed with 75% ethanol and dried at room temperature. LCM was performed on the same day using a laser dissection microscope (Leica, LMD 7000) with the following settings: excitation wavelength 495 nm, laser power 30, aperture 5, speed 6 and pulse frequency 119. From each animal, at least 800 plaque cores and halos, dissected from 12 brain sections were cut using a 63 \times magnification objective, collected in 0.5 ml caps (Leica Microsystem) and subsequently pooled for protein analysis. Areas containing no plaques were cut using a 10 \times magnification objective and were used as controls. Protein lysates were done essentially as described above using RIPA with 0.1% SDS.

Slice preparation and electrophysiological recordings applying A η peptides *in vitro*. Transverse hippocampal slices (350 μm) were prepared from P20–30 Swiss mice following standard procedures⁵². Slices were cut in ice-cold oxygenated (95% O₂, 5% CO₂) solution containing 206 mM sucrose, 2.8 mM KCl, 1.25 mM NaH₂PO₄, 2 mM MgSO₄, 1 mM MgCl₂, 1 mM CaCl₂, 26 mM NaHCO₃, 0.4 mM sodium ascorbate and 10 mM glucose, pH 7.4. For recovery (1 h), slices were incubated at 27°C in oxygenated standard ACSF containing: 124 mM NaCl, 2.8 mM KCl, 1.25 mM NaH₂PO₄, 2 mM MgSO₄, 3.6 mM CaCl₂, 26 mM NaHCO₃, 0.4 mM sodium ascorbate, 10 mM glucose (pH 7.4)⁵³. Slices were inspected in a chamber on an upright microscope (Slicescope, Scientifica Ltd) with infrared differential interference contrast illumination, and were perfused with the oxygenated ACSF at $27 \pm 1^\circ\text{C}$. fEPSPs were recorded in the stratum radiatum of the CA1 region using a glass electrode (filled with 1 M NaCl, 10 mM HEPES, pH 7.4) and the stimuli (30% of maximal fEPSP) were delivered to the Schaeffer Collateral pathway by a monopolar glass electrode (filled with ACSF). Electrodes were specifically placed just below the surface of the slice to maximize the exposure to circulating peptides. A minimum of 15–20 min stable baseline was first obtained in standard ACSF followed by another 15–20 min of bath application of ACSF containing SEC fractions (CHO, A η - α or A η - β ; 1/15 dilution, interleaved recordings) using re-circulation with a peristaltic pump at $2.5\text{--}3\text{ ml min}^{-1}$ while being continuously aerated with 95% oxygen. No alterations in fEPSP baseline responses were observed after incubation with the SEC fractions (Extended Data Fig. 8). In the continuous presence of ACSF/SEC solution, LTP was induced using a high-frequency stimulation protocol with two pulses of 100 Hz for 1 s with a 20 s interval between pulses, and recorded for 1 h. Control recordings (no application of SEC fractions) were obtained in an interleaved fashion in which ACSF was re-circulated using an identical procedure. For LTP analysis, the first third of the fEPSP slope was calculated in baseline condition (15–20 min before induction of LTP) and compared to that after LTP induction (60 min after tetanization of Schaeffer collaterals). The average baseline value was normalized to 100% and all values of the experiment were normalized to this baseline average (1-min bins). Experimental data were pooled per condition and presented as mean \pm s.e.m. Data analysis was performed with the Clampfit software (Molecular Devices). The test samples were blinded to the investigator and uncoded at the end of the experiments. Statistical analysis was performed using GraphPad (Prism 6) with the last 15 min of the recordings compared to measurements of 15–20 min of baseline, using a two-tailed Student's *t*-test for statistical analysis on two samples or one-way ANOVA and post hoc Bonferroni test for statistical analysis on three and more samples, with $P < 0.05$ taken as statistically significant. No power analysis was done to estimate sample size, and there was no randomization.

Electrophysiological recordings of the effects of BACE1 inhibitor *in vitro*. BACE1 inhibitor (100 mg kg⁻¹, single gavage) or vehicle-treated mice (12 weeks old) were deeply anaesthetized with isoflurane (1% in O₂) and decapitated with brains rapidly extracted and placed for 5–6 min in ice-cold bubbled (95% O₂, 5% CO₂) slicing solution (in mM): 75 sucrose, 85 NaCl, 2.5 KCl, 1.25 NaH₂PO₄, 25 NaHCO₃, 0.5 CaCl₂, 4 MgCl₂, 25 glucose, pH 7.4. Coronal slices (400 μm) containing the hippocampus were cut (VT1200S; Leica) and transferred into a warming chamber (35°C) filled with bubbled solution of the same composition, except sucrose was omitted and NaCl increased to 125 mM (30 min). This was followed by the transfer of slices into recording ACSF (in mM): 125 NaCl, 2.5 KCl, 1.25 NaH₂PO₄, 25 NaHCO₃, 2 CaCl₂, 2 MgCl₂, 25 glucose. Recordings of fEPSP were made from the hippocampal CA1 area. A glass bipolar stimulating electrode was placed in the stratum radiatum of CA2–CA3 subfields to stimulate Schaeffer collaterals with 0.2 ms current pulses at 0.033 Hz (A-360, WPI), with evoked responses recorded in the stratum radiatum of CA1 area. Incrementing current pulses (0.2 mA) were used for obtaining stimulus-response relationship graphs. Stable baseline and LTP recordings were made using one-half of the maximal stimulus intensities; LTP was induced by high-frequency stimulation of Schaeffer collaterals with 10 trains of 10 pulses at 100 Hz applied, with 2-s inter-train intervals. Signals were filtered at 5 kHz, digitally sampled at 10 kHz and stored for offline analysis. The relative slope and peak amplitude of evoked fEPSPs were measured using FitMaster (HEKA Electronics). The groups of treated mice were blinded to the investigator and uncoded at the end of the experiments. A one-way ANOVA and post hoc Bonferroni test have been used for statistical analysis, with $P < 0.05$ taken as statistically significant.

***In vivo* two-photon Ca²⁺ imaging.** All experimental procedures were in compliance with institutional animal welfare guidelines and were approved by the state government of Bavaria, Germany. The animal preparation procedure was similar to that described previously²⁵. In brief, C57Bl/6 mice (male or female, \sim P40) were anaesthetized with isoflurane (1–1.5%) and placed onto a warming plate ($37\text{--}38^\circ\text{C}$). The skin was removed and a custom-made recording chamber was glued to the exposed skull. A craniotomy (\sim 1 mm) was made over the

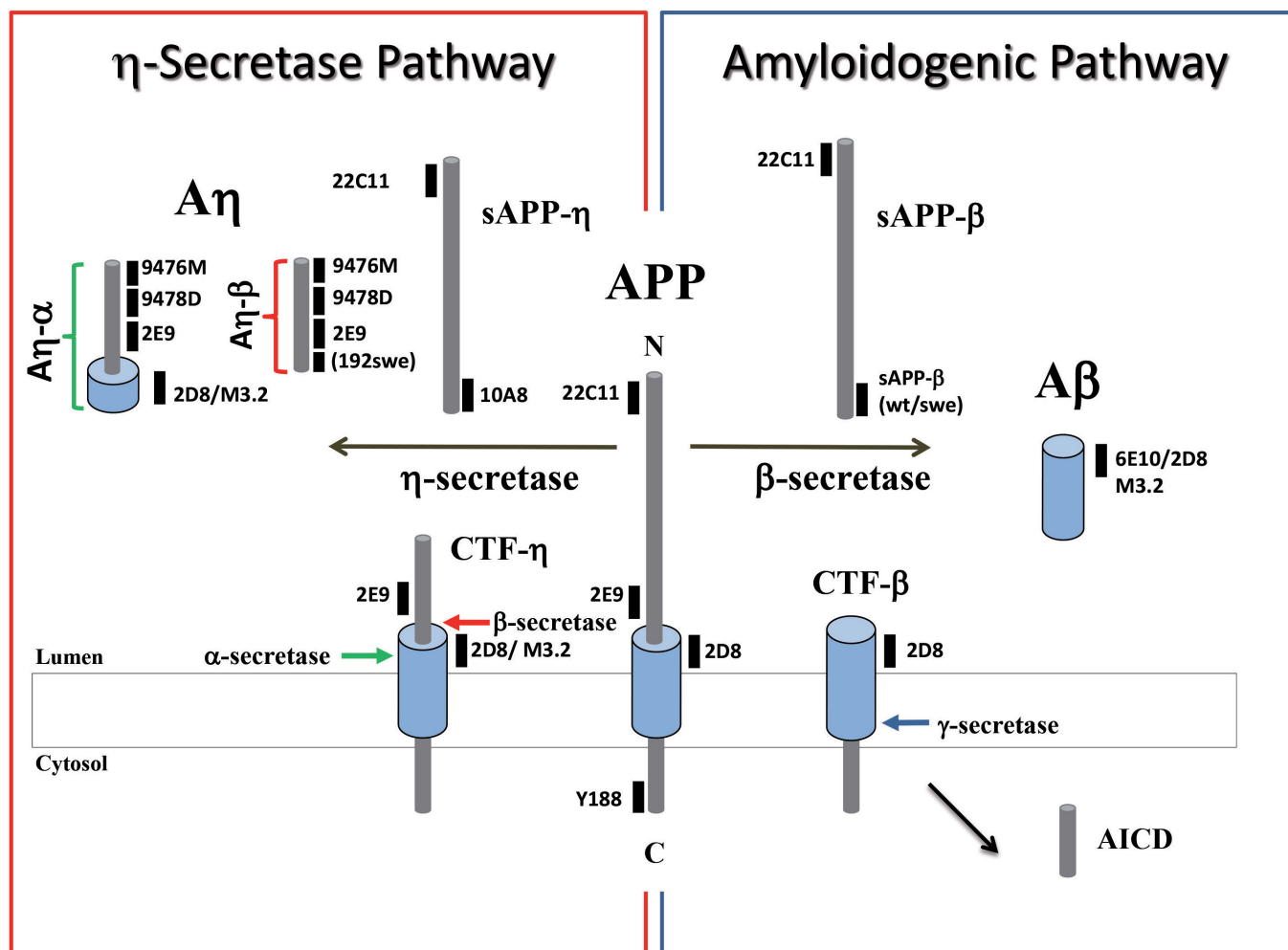
hippocampus (2.5 mm posterior to bregma, 2.2 mm lateral to the midline) and a small portion of the overlying cortical tissue was carefully removed by aspiration. The animal was placed under a microscope on a warm heating plate (37–38 °C) and kept anaesthetized with low-levels of isoflurane (~0.8%). Respiratory and pulse rates were continuously monitored. The recording chamber was perfused with warm normal Ringer's solution containing 125 mM NaCl, 4.5 mM KCl, 26 mM NaHCO₃, 1.25 mM NaH₂PO₄, 2 mM CaCl₂, 1 mM MgCl₂ and 20 mM glucose (pH 7.4 when bubbled with 95% O₂ and 5% CO₂). The exposed CA1 region of the hippocampus was then stained with fluo-8AM (AAT Bioquest; 0.6 mM) using the multi-cell bolus loading technique⁵⁴.

In vivo imaging was performed with a custom-built two-photon microscope equipped with a Ti:sapphire laser system (Coherent; laser wavelength 925 nm), a resonant scanner and a Pockel's cell for laser intensity modulation. Full-frame images were acquired at 30 Hz using a water-immersion objective (Nikon; 40×, 0.8 numerical aperture). Data acquisition was controlled using custom-written software based on LabVIEW (National Instruments). Image analysis was performed off-line by using custom routines in LabVIEW and Igor Pro (Wavemetrics). Cellular regions of interest were drawn around individual somata, and then relative fluorescence change ($\Delta F/F$) versus time traces were generated for each region of interest. Ca²⁺ transients were identified as changes in $\Delta F/F$ that were three times larger than the s.d. of the noise band.

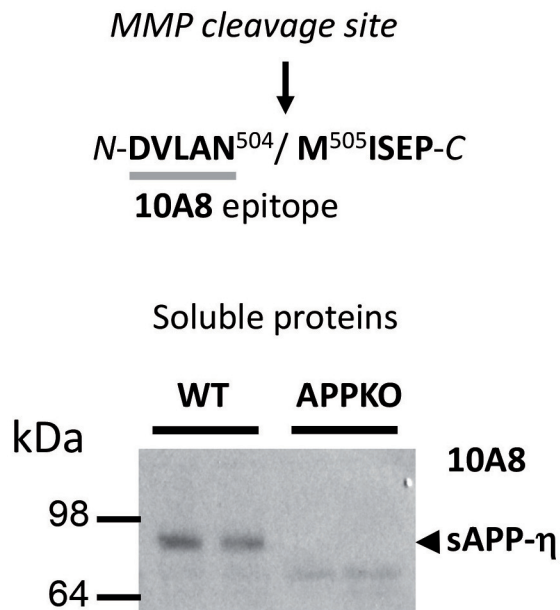
To assess the effects of A η peptides on neuronal activity *in vivo*, the peptides or the respective controls (SEC fractions obtained from untransfected CHO cells or a synthetic peptide (46 amino acids; 1-ADSVPANTENEVEPVDARPAADRGLTTRPGSGLTNIKTEEISEVKM-46) of a middle part of A η were added to the normal Ringer's solution used for perfusion of the recording chamber (bath-application technique; 45–60 min each wash-in). In a subset of experiments, synthetic A η - α (92 amino acids; 1-MISEPRISYGNDALMPSLTETKTTVELLPVNGEFLDDLQPWHSFGADSVANTENEVEPVDARPAADRGLTTRPGSGLTNIKTEEISEVKM-92) was applied locally by gentle pressure injection through a glass pipette that was placed close to the neurons of interest (local application technique; 40 s each pressure injection).

The samples were blinded to the investigator and uncoded at the end of the experiments. Statistical analysis was performed using SPSS. The statistical methods used were the Student's *t*-test and the Fisher's exact test. *P* < 0.05 was considered statistically significant.

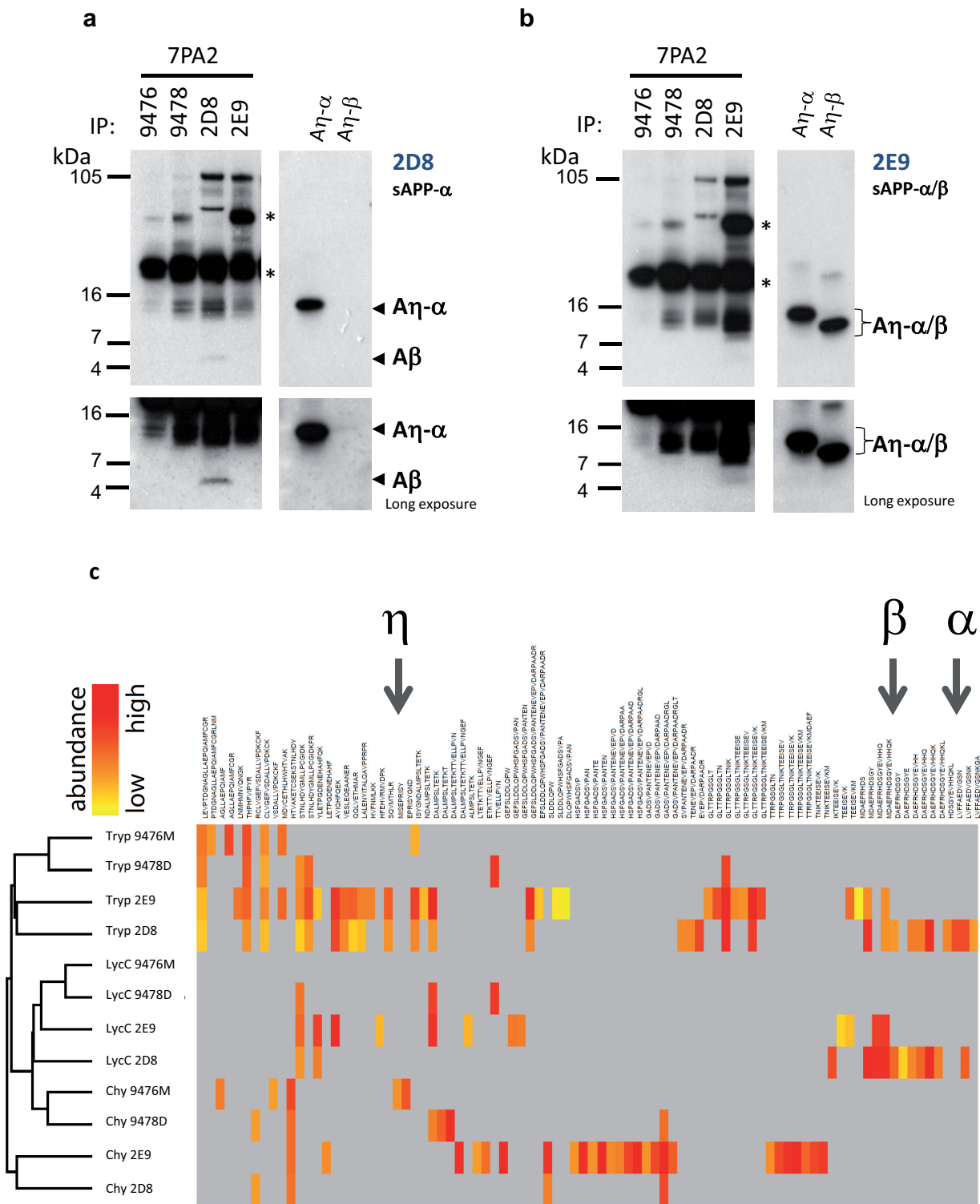
31. Podlisny, M. B. *et al.* Aggregation of secreted amyloid β -protein into sodium dodecyl sulfate-stable oligomers in cell-culture. *J. Biol. Chem.* **270**, 9564–9570 (1995).
32. Kaech, S. & Banker, G. Culturing hippocampal neurons. *Nature Protocols* **1**, 2406–2415 (2006).
33. Israel, M. A. *et al.* Probing sporadic and familial Alzheimer's disease using induced pluripotent stem cells. *Nature* **482**, 216–220 (2012).
34. Shi, Y., Kirwan, P., Smith, J., Robinson, H. P. & Livesey, F. J. Human cerebral cortex development from pluripotent stem cells to functional excitatory synapses. *Nature Neurosci.* **15**, 477–486 (2012).
35. Jacobsen, H. *et al.* Combined treatment with a BACE inhibitor and anti-Ab antibody gantenerumab enhances amyloid reduction in APP^{London} mice. *J. Neurosci.* **34**, 11621–11630 (2014).
36. Nolan, R. L. & Teller, J. K. Diethylamine extraction of proteins and peptides isolated with a mono-phasic solution of phenol and guanidine isothiocyanate. *J. Biochem. Biophys. Methods* **68**, 127–131 (2006).
37. Westmeyer, G. G. *et al.* Dimerization of β -site β -amyloid precursor protein-cleaving enzyme. *J. Biol. Chem.* **279**, 53205–53212 (2004).
38. Rappsilber, J., Mann, M. & Ishihama, Y. Protocol for micro-purification, enrichment, pre-fractionation and storage of peptides for proteomics using StageTips. *Nature Protocols* **2**, 1896–1906 (2007).
39. Scheltema, R. A. *et al.* The Q Exactive HF, a Benchtop mass spectrometer with a pre-filter, high-performance quadrupole and an ultra-high-field Orbitrap analyzer. *Mol. Cell. Proteomics* **13**, 3698–3708 (2014).
40. Olsen, J. V. *et al.* Higher-energy C-trap dissociation for peptide modification analysis. *Nature Methods* **4**, 709–712 (2007).
41. Cox, J. & Mann, M. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nature Biotechnol.* **26**, 1367–1372 (2008).
42. Cox, J. *et al.* Andromeda: a peptide search engine integrated into the MaxQuant environment. *J. Proteome Res.* **10**, 1794–1805 (2011).
43. R Development Core Team. *R: A Language and Environment for Statistical Computing* (R Foundation for Statistical Computing, <http://www.R-project.org/> (2014)).
44. Haass, C. *et al.* The Swedish mutation causes early-onset Alzheimer's-disease by β -secretase cleavage within the secretory pathway. *Nature Med.* **1**, 1291–1296 (1995).
45. McKhann, G. *et al.* Clinical diagnosis of Alzheimer's disease: report of the NINCDS-ADRDA Work Group under the auspices of Department of Health and Human Services Task Force on Alzheimer's Disease. *Neurology* **34**, 939–944 (1984).
46. Lannfelt, L. *et al.* Amyloid precursor protein mutation causes Alzheimer's disease in a Swedish family. *Neurosci. Lett.* **168**, 254–256 (1994).
47. Lannfelt, L. *et al.* Amyloid β -peptide in cerebrospinal fluid in individuals with the Swedish Alzheimer amyloid precursor protein mutation. *Neurosci. Lett.* **199**, 203–206 (1995).
48. Braak, H. & Braak, E. Neuropathological staging of Alzheimer-related changes. *Acta Neuropathol.* **82**, 239–259 (1991).
49. Mirra, S. S. *et al.* The Consortium to Establish a Registry for Alzheimer's Disease (CERAD). Part II. Standardization of the neuropathologic assessment of Alzheimer's disease. *Neurology* **41**, 479–486 (1991).
50. Hyman, B. T. *et al.* National Institute on Aging-Alzheimer's Association guidelines for the neuropathologic assessment of Alzheimer's disease. *Alzheimers Dement.* **8**, 1–13 (2012).
51. Liao, L. *et al.* Proteomic characterization of postmortem amyloid plaques isolated by laser capture microdissection. *J. Biol. Chem.* **279**, 37061–37068 (2004).
52. Houeland, G. *et al.* Transgenic mice with chronic NGF deprivation and Alzheimer's disease-like pathology display hippocampal region-specific impairments in short- and long-term plasticities. *J. Neurosci.* **30**, 13089–13094 (2010).
53. Townsend, M., Shankar, G. M., Mehta, T., Walsh, D. M. & Selkoe, D. J. Effects of secreted oligomers of amyloid beta-protein on hippocampal synaptic plasticity: a potent role for trimers. *J. Physiol. (Lond.)* **572**, 477–492 (2006).
54. Stosiek, C., Garaschuk, O., Holthoff, K. & Konnerth, A. *In vivo* two-photon calcium imaging of neuronal networks. *Proc. Natl Acad. Sci. USA* **100**, 7319–7324 (2003).



Extended Data Figure 1 | Schematic presentation of the η -secretase processing pathway. Schematic representation of the η -secretase pathway (left) as compared to the amyloidogenic pathway (right). Antibodies used in this study are indicated.

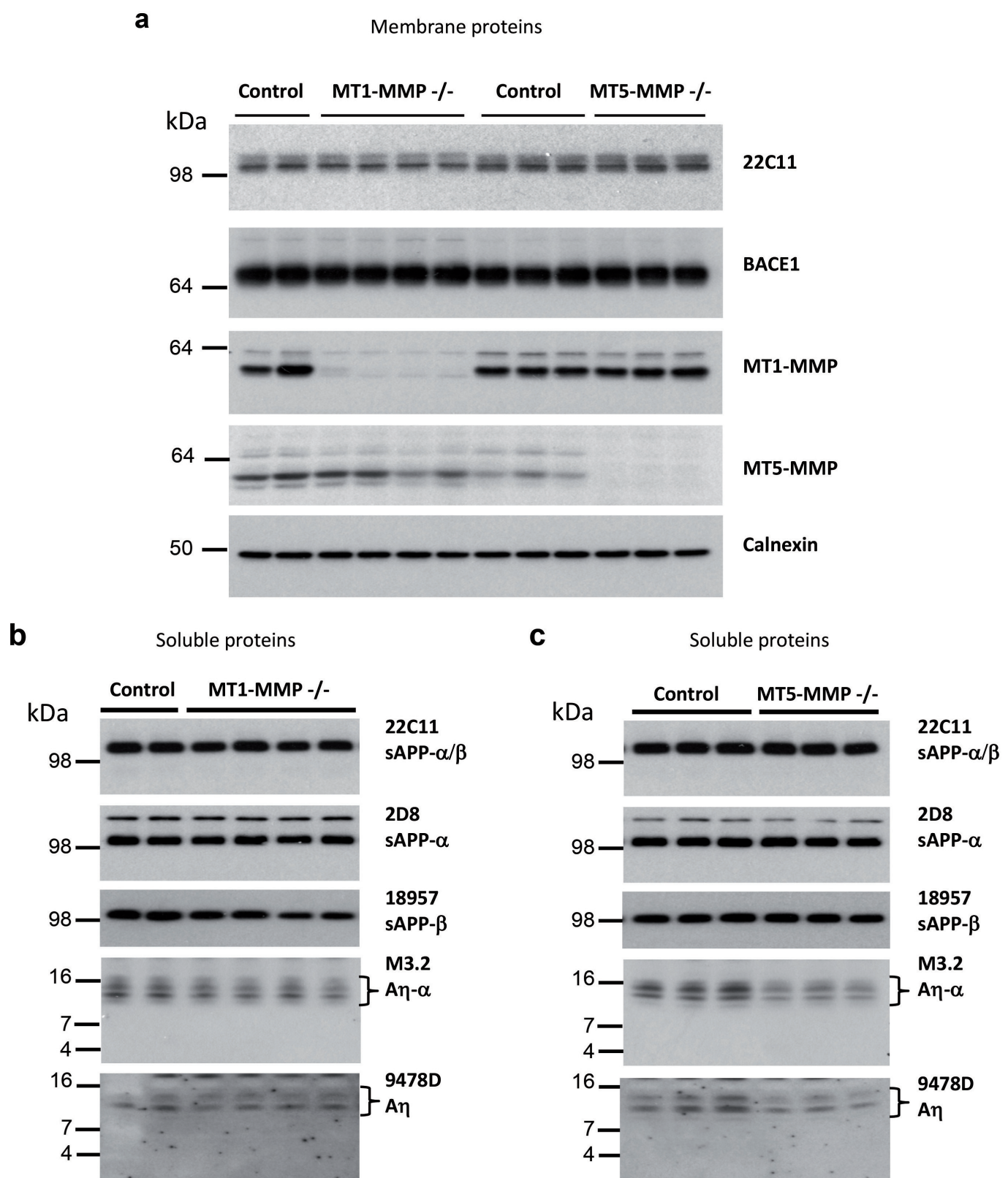


Extended Data Figure 2 | η-Secretase cleavage at Met505 of APP₆₉₅. MMP proteins can cleave human APP₆₉₅ at the indicated position (arrow) between amino acids N504 and M505 in the N-terminal domain. The epitope for the neo-epitope-specific antibody 10A8 is indicated (grey line). sAPP-η was specifically detected in diethylamine (DEA; 0.2% diethylamine in 50 mM NaCl, pH 10) extracts of P10 wild-type mouse brain using antibody 10A8, but was absent in APP-knockout brains. Of note, antibody 10A8 failed to detect sAPP-α/β, confirming its selectivity for the η-cleavage site.



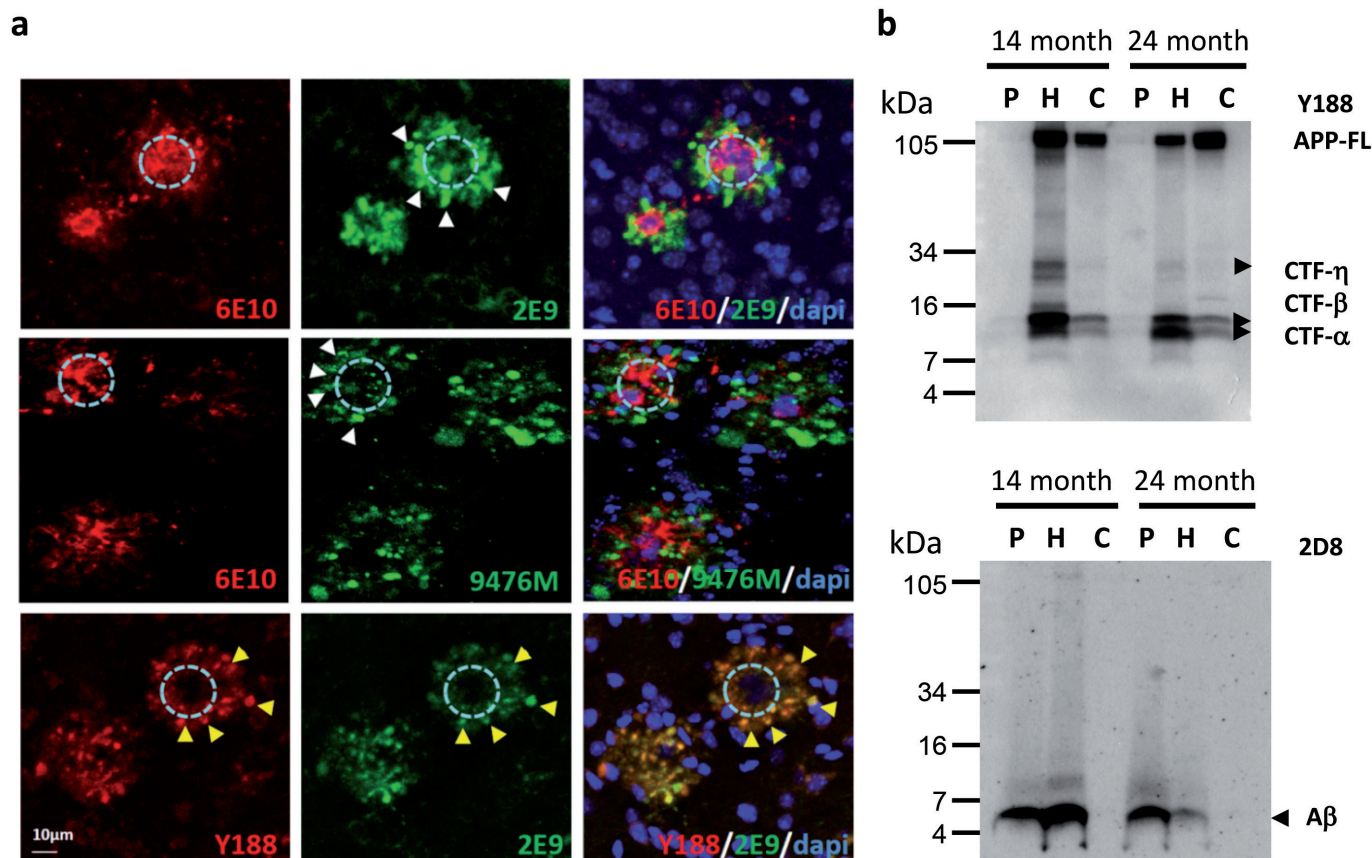
a, b, After removal of sAPP- α from conditioned media of CHO 7PA2 cells using appropriate centricon filters, the flow-through was used to isolate A η peptides by immunoprecipitation. Synthetic peptides (1 ng per lane) were loaded to indicate the respective sizes of A η - α and A η - β . A η peptides were captured with antibodies 9476M and 9478D directed against the putative N-terminal epitopes (Supplementary Table 1), 2E9 against a middle domain of A η and 2D8 (which also immunoprecipitates amyloid- β). 2D8 detection (**a**) revealed that all antibodies captured peptides positive for the N-terminal part of the amyloid- β domain in a molecular mass range of synthetic A η - α

between 12 and 16 kDa. The same samples analysed with 2E9 (**b**) confirmed the presence of A η in all samples, with the lowest levels when precipitated with 9476M. Note that unlike 2D8 antibody, 2E9 allows the additional detection of A η - β . (Asterisks denote IgG.) **c**, A heat map of peptides identified after analytic proteolysis by mass spectrometry analysis, with 7PA2 supernatants immunoprecipitated with 2D8, 2E9, 9476M and 9478D antibodies. Arrows indicate peptides that start exactly with the amino acid sequence C-terminal of the respective cleavage sites of the η -secretase, β -secretase or α -secretase site. Chy, chymotrypsin; LysC, protease LysC; tryp, trypsin.



Extended Data Figure 4 | MT5-MMP displays η -secretase activity in brain. **a–c**, *MT1-MMP* $^{-/-}$ and *MT5-MMP* $^{-/-}$ (also known as *Mmp14* $^{-/-}$ and *Mmp24* $^{-/-}$, respectively) mice were analysed for changes in η -secretase activity. Membrane and soluble proteins from P10 mouse brains were analysed. **a**, In RIPA lysates, no changes in APP and BACE1 levels were detected in knockout brains. MT1- and MT5-MMP were selectively knocked out as shown

by the lack of signals in western blots. Calnexin served as a loading control. Soluble A η levels, detected by antibodies 9478D and M3.2 (A η - α) were unchanged in *MT1-MMP* $^{-/-}$ mouse brains (**b**), but reduced in *MT5-MMP* $^{-/-}$ mouse brains (**c**). Total levels of secreted APP (22C11), sAPP- α or sAPP- β were unchanged (**b**, **c**).

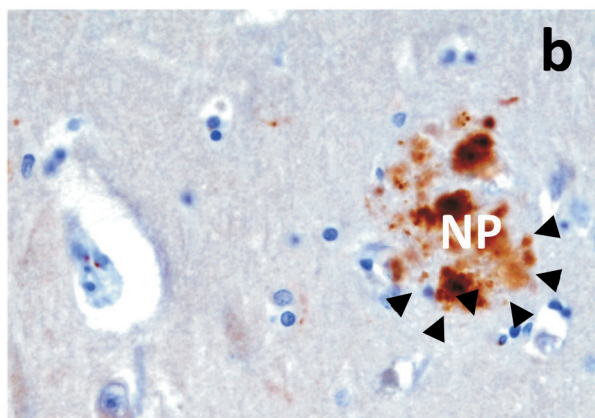
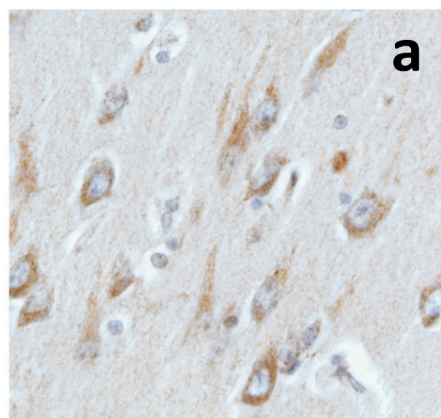


Extended Data Figure 5 | Accumulation of CTF- η in dystrophic neurites.
a, Immunohistological stainings of cortical sections of 6-month-old APPPS1-21 transgenic mice ($n = 3$) revealed 6E10-positive amyloid- β plaque cores (encircled) surrounded by dystrophic neurites positive for 2E9 (white arrowheads, top) and 9476M (white arrowheads, middle). Y188 (bottom panel) staining co-localized with 2E9-positive signal (yellow arrowheads, bottom). Nuclei were counterstained with DAPI. Scale bar, 10 μ m. **b**, Accumulation of CTF- η fragment in dystrophic neurites of 14-month and 24-month-old APPPS1-21 mice. Western blot analysis of material obtained by LCM of

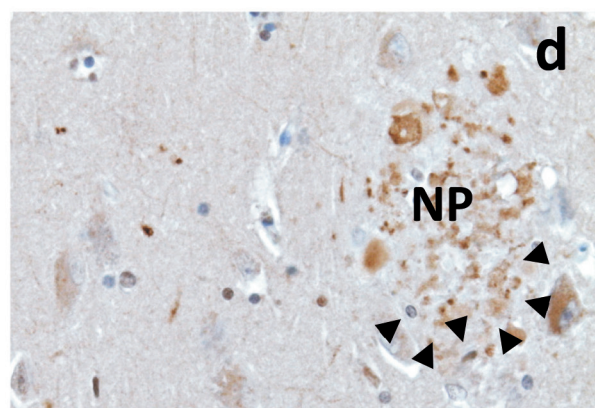
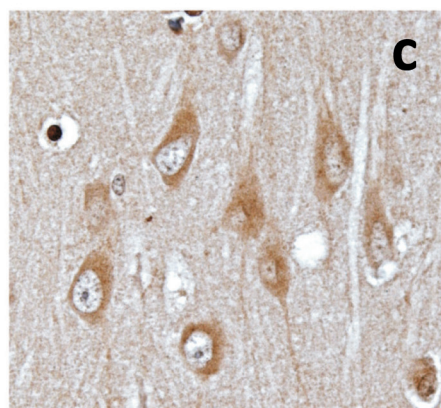
APPPS1-21 brain sections ($n = 5$) bearing thioflavin-S-positive amyloid- β plaque core (P) and the surrounding amyloid- β plaque halo (H). As a control (C), brain areas devoid of plaques were used. While amyloid- β was readily detected by antibody 2D8 in lysates containing plaque-enriched material and halo regions (fractions P and H; bottom), CTF- η was selectively detected in the lysates prepared from the region enriched in dystrophic neurites (H; top), but not detected in plaque or control regions (P or C). As expected, CTF- β / α species are also enriched in dystrophic neurites (H).

Control

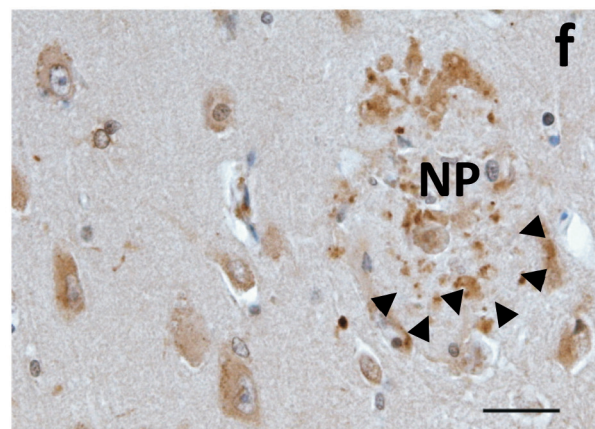
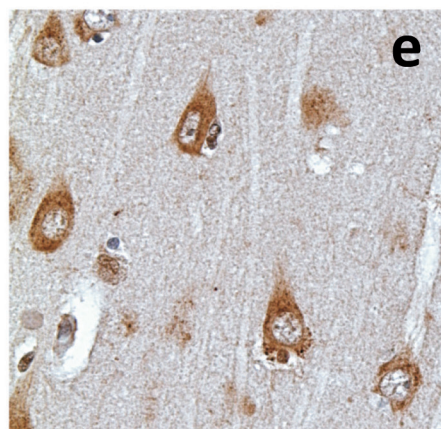
AD



22C11



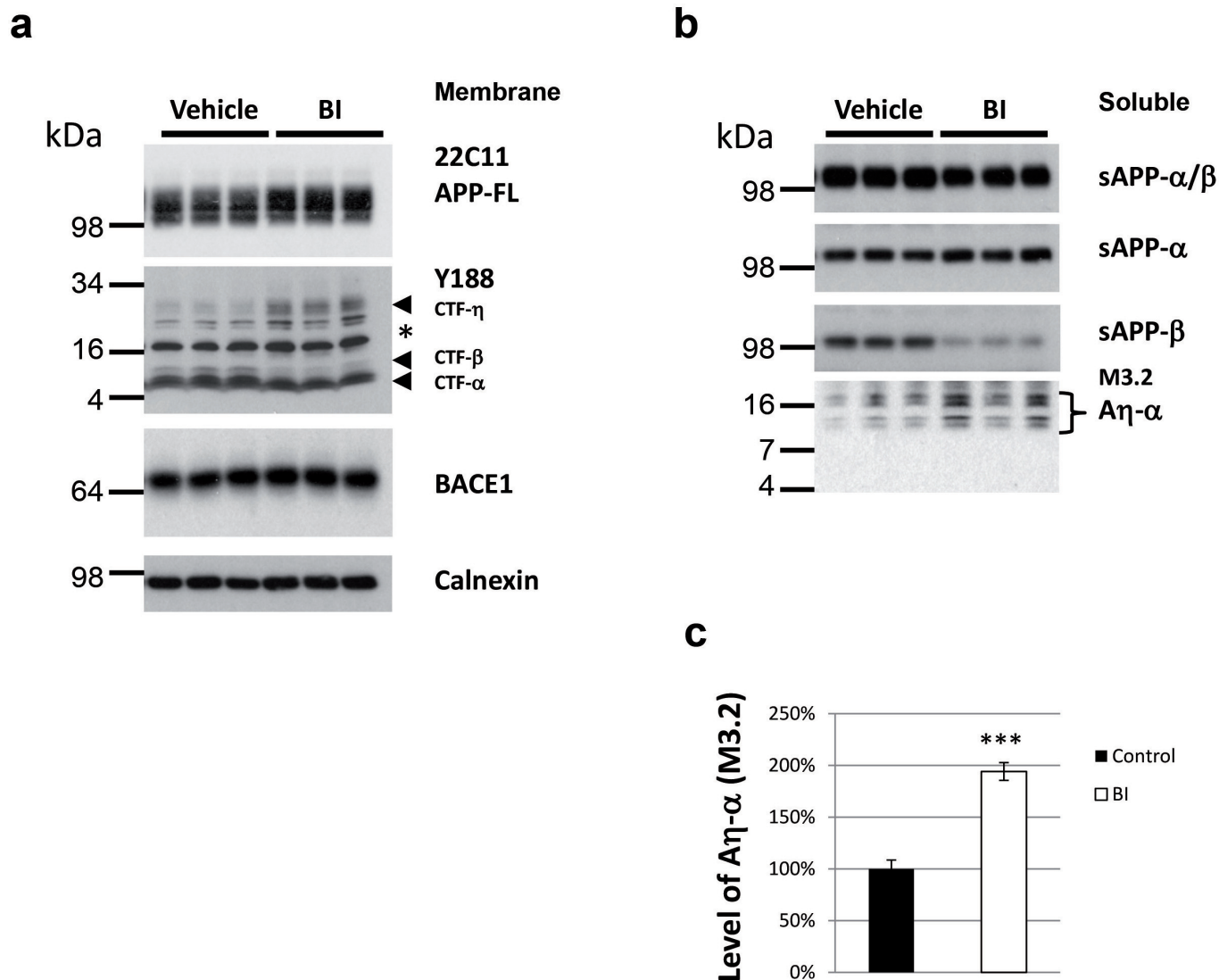
9478D



9476M

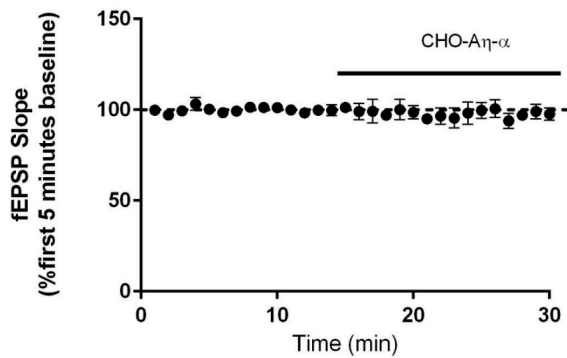
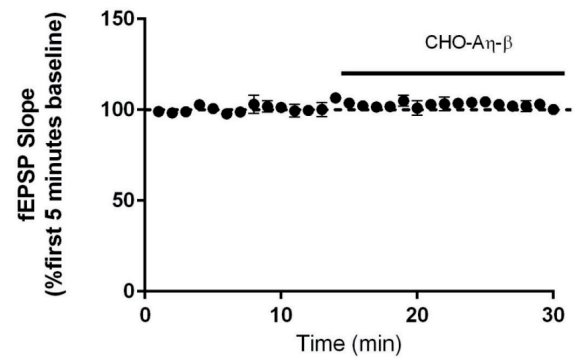
Extended Data Figure 6 | Dystrophic neurites in AD brains are positive for A η -epitope antibodies. a–f, Immunohistochemistry with 22C11 (a, b), 9478D (c, d) and 9476M (e, f) antibodies in the human hippocampus (CA1-subiculum region) of a control case (a, c, e) and an AD case (b, d, f).

Immuno-positive signals were observed with 22C11 (a), 9478D (c) and 9476M (e) antibodies in the somata and neuropils of a normal and AD brain. In AD brains, these antibodies decorate dystrophic neurites (b, d, f, denoted by arrowheads). Scale bar, 30 μ m. NP, neuritic plaque.



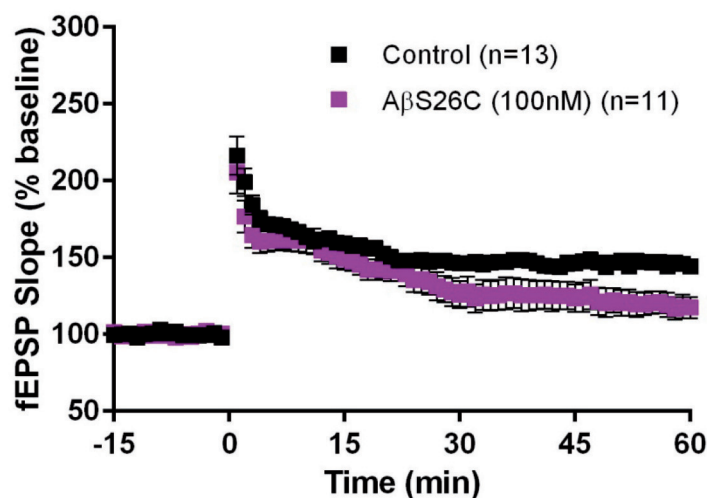
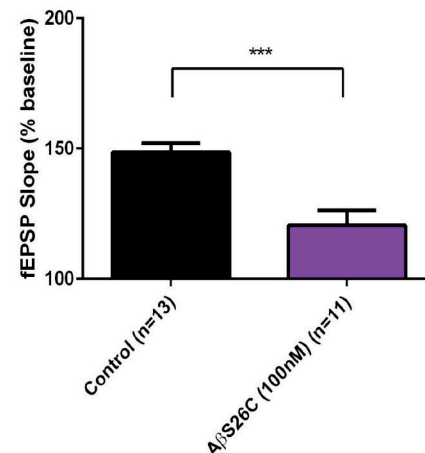
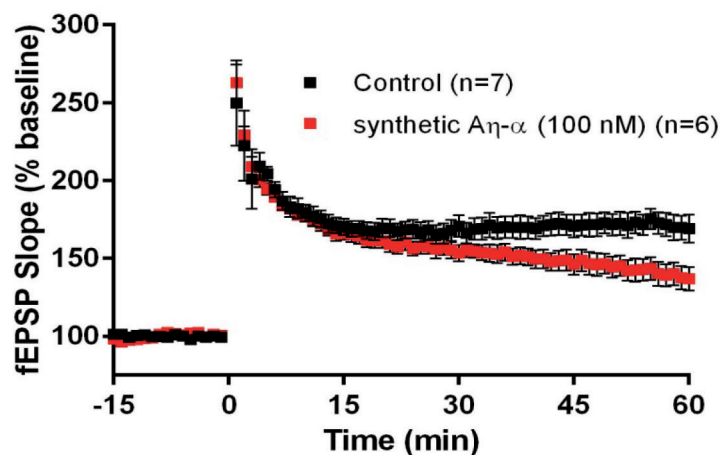
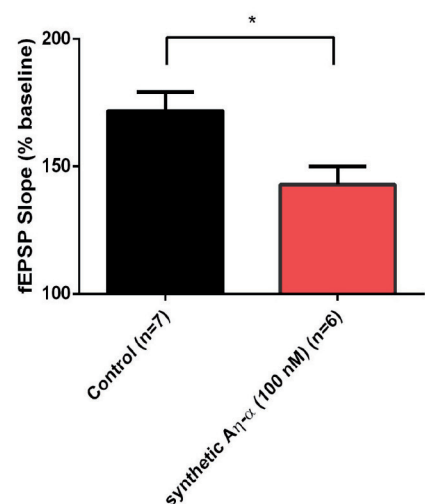
Extended Data Figure 7 | Increased A η levels after acute treatment with a BACE1 inhibitor. **a**, In membrane lysates of brains obtained from animals treated with BACE1 inhibitor (BI, 100 mg kg⁻¹ SCH1682496), an increase in CTF- η was observed, which was paralleled by a strong reduction of CTF- β , while CTF- α was unchanged. APP-FL and BACE1 signals remained

unchanged (asterisk indicates background band). Calnexin served as a loading control. **b**, In the soluble fraction, BACE1 inhibition resulted in enhanced A η - α levels and reduced sAPP- β levels, indicating efficient BACE1 inhibition. **c**, Production of A η - α species, which was detected by antibody M3.2, revealed a 95.4% increase after BACE1 inhibition. $n = 3$; $P < 0.01$, Student's t -test.

a**b**

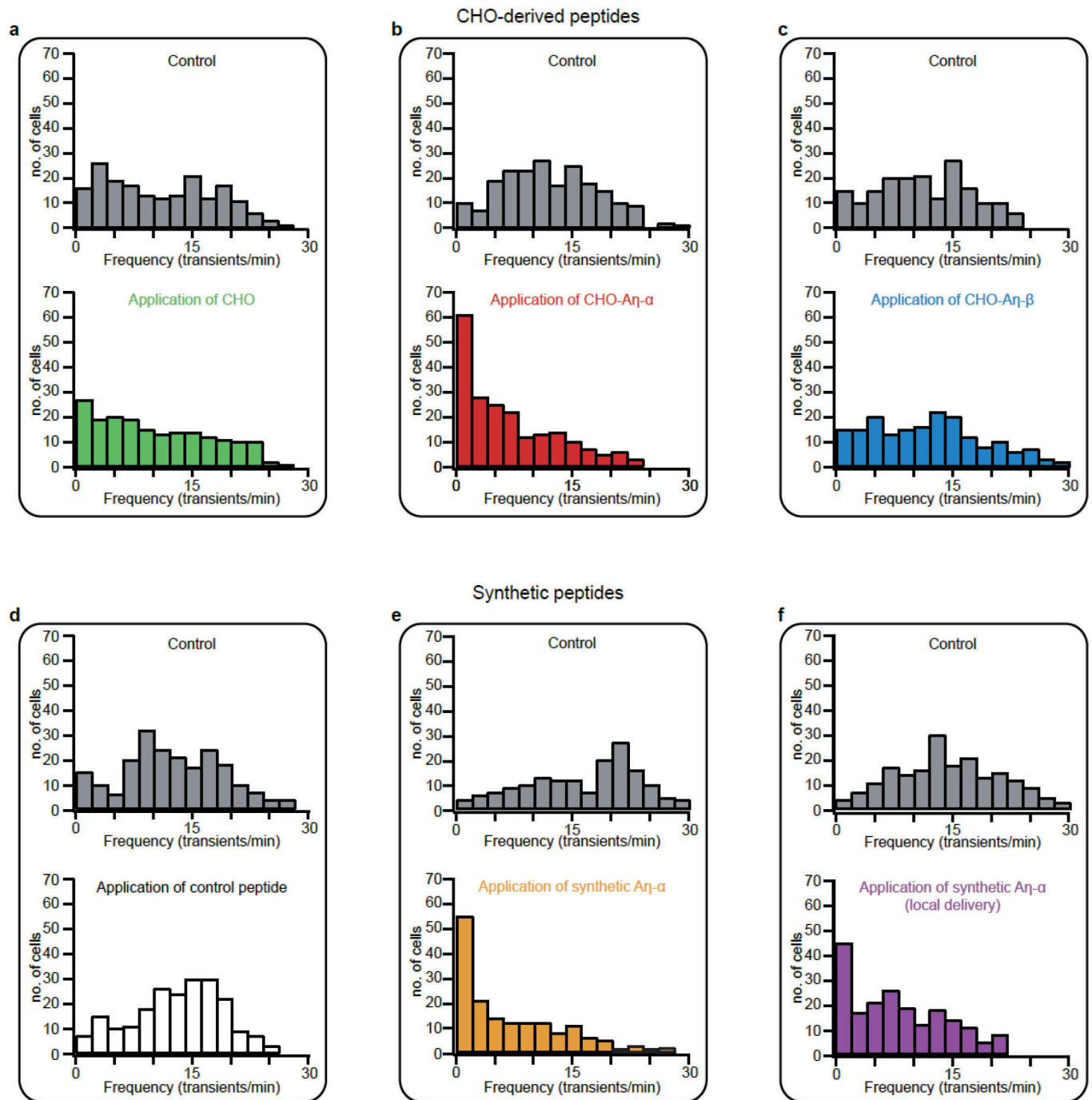
Extended Data Figure 8 | A η - α and A η - β derived from CHO cells did not influence baseline activity at the hippocampal CA3–CA1 synapse. Soluble A η - α and A η - β peptides were expressed in CHO cells and collected in OPTIMEM medium. **a**, **b**, SEC fractions containing A η were diluted (1:15) in ACSF for the treatment of hippocampal slices and LTP measurements. A η - α or

A η - β SEC fractions were perfused over mouse hippocampal slices after obtaining a 15-min stable baseline of a fEPSP at the CA3–CA1 synapse. The baseline remained unchanged for another 15 min when slices were incubated with CHO-cell-derived recombinant A η - α (**a**) or A η - β (**b**).

a**b****c****d**

Extended Data Figure 9 | Aβ_{S26C} dimers and synthetic Aη-α impair hippocampal LTP. **a**, In line with previous findings²³, Aβ_{S26C} cross-linked dimers (containing cysteine instead of serine at residue 26; 100 nM final; JPT Peptide Technologies; diluted in 25 ml re-circulating ACSF) reduced LTP as compared to interleaved control LTP recordings in 25 ml re-circulating ACSF. **b**, Illustrated is the average LTP magnitude (at 45–60 min after LTP induction) normalized to pre-LTP baseline values (100%) in untreated and

treated conditions (*** $P < 0.001$; Student's t -test). **c**, Treatment with synthetic Aη-α (100 nM final; Peptide Speciality Laboratories; diluted in 25 ml re-circulating ACSF) reduced LTP as compared to interleaved control LTP recordings in 25 ml re-circulating ACSF. **d**, Illustrated is the average LTP magnitude (at 45–60 min post-LTP induction) normalized to pre-LTP baseline values (100%) in treated and untreated conditions (* $P < 0.05$; Student's t -test).



Extended Data Figure 10 | A η - α decreases the frequencies of neuronal calcium transients *in vivo*. a–f, Histograms showing in each panel the corresponding distributions of calcium transients before (control) and during subsequent exposure of A η peptides (b, c, e, f) and controls (a, d).

Inhibition of Gli1 mobilizes endogenous neural stem cells for remyelination

Jayshree Samanta¹, Ethan M. Grund¹, Hernandez M. Silva², Juan J. Lafaille², Gord Fishell¹ & James L. Salzer¹

Enhancing repair of myelin is an important but still elusive therapeutic goal in many neurological disorders¹. In multiple sclerosis, an inflammatory demyelinating disease, endogenous remyelination does occur but is frequently insufficient to restore function. Both parenchymal oligodendrocyte progenitor cells and endogenous adult neural stem cells resident within the subventricular zone are known sources of remyelinating cells². Here we characterize the contribution to remyelination of a subset of adult neural stem cells, identified by their expression of Gli1, a transcriptional effector of the sonic hedgehog pathway. We show that these cells are recruited from the subventricular zone to populate demyelinated lesions in the forebrain but never enter healthy, white matter tracts. Unexpectedly, recruitment of this pool of neural stem cells, and their differentiation into oligodendrocytes, is significantly enhanced by genetic or pharmacological inhibition of Gli1. Importantly, complete inhibition of canonical hedgehog signalling was ineffective, indicating that the role of Gli1 both in augmenting hedgehog signalling and in retarding myelination is specialized. Indeed, inhibition of Gli1 improves the functional outcome in a relapsing/remitting model of experimental autoimmune encephalomyelitis and is neuroprotective. Thus, endogenous neural stem cells can be mobilized for the repair of demyelinated lesions by inhibiting Gli1, identifying a new therapeutic avenue for the treatment of demyelinating disorders.

Remyelination in the adult human and mouse brains is performed by two cell types: oligodendrocyte progenitor cells (OPCs) and neural stem cells (NSCs). OPCs, which are present in the parenchyma of healthy brain as well as in, and around, multiple sclerosis lesions³, can be identified by their expression of the NG2 proteoglycan and platelet-derived growth factor receptor alpha (PDGFR- α). They respond locally to demyelination by generating oligodendrocytes although do not migrate long distances during remyelination^{4,5}. NSCs present in the subventricular zone (SVZ), express glial fibrillary acidic protein (GFAP) and Nestin, and are normally quiescent. In response to demyelination, cells in the adult SVZ can generate oligodendrocytes⁶, including, presumptively, in patients with multiple sclerosis⁷.

The signals that activate and recruit NSCs to lesion sites and promote their local differentiation into oligodendrocytes remain poorly understood. A candidate to regulate NSCs is sonic hedgehog (Shh), an important morphogen during CNS development that is required for the generation of most oligodendrocytes during development⁸ and for the maintenance of stem cells in the adult SVZ⁹. Shh is therefore an attractive candidate to expand the pool of premyelinating cells available for repair. Indeed, Shh levels have been reported to increase in remyelinating lesions¹⁰.

Canonical Shh signalling is mediated by interactions of the hedgehog receptor patched (Ptc) with the G-protein-coupled transmembrane co-receptor smoothened (Smo). Binding of Shh to Ptc relieves its inhibition of Smo and thereby activates the Gli family of zinc-finger transcription factors¹¹. Of the three Gli proteins, Gli1 is the only one

whose transcription is driven by Shh signalling and its expression is therefore considered a sensitive readout of sustained, high-level activation of this pathway^{12,13}.

In this study, we have examined remyelination by the Shh-responsive (that is, Gli1⁺) pool of NSCs, which is concentrated in the ventral SVZ and comprises ~25% of NSCs¹⁴. To genetically fate-map Gli1⁺ NSCs, we crossed *Gli1*^{CreERT2/+} mice with the *Rosa-CAG-EGFP* (RCE) reporter¹⁵ to generate *Gli1*^{CE} mice in which tamoxifen treatment results in permanent expression of cytoplasmic green fluorescent protein (GFP) in all Gli1-expressing cells and their progeny^{12,14}; see Supplementary Table 1 for a summary of these and other mouse lines used in this study. GFP-labelled cells correspond to NSCs in the SVZ and a subset of astrocytes, but not OPCs or oligodendrocytes. We then followed the fate of these GFP⁺ cells after inducing demyelination in the mouse corpus callosum (CC) either by (1) dietary cuprizone¹⁶ or (2) direct, stereotactic injection of the detergent lysophosphatidyl-choline (LPC)⁴.

At 6 weeks of dietary cuprizone, corresponding to peak demyelination, GFP-expressing cells were recruited to areas of demyelination. In contrast, no labelled cells were observed in the CC of controls (Fig. 1a). At 2 weeks of remyelination, after removal of cuprizone from the diet, GFP-expressing cells in the CC (17.3 cells \pm 2.6 per section) differentiated exclusively into glia, primarily oligodendroglia, namely PDGFR- α ⁺ OPCs (9.8 \pm 8.7%) and CC1⁺ oligodendrocytes (40.2 \pm 15.1%), as well as GFAP-expressing astrocytes (15.5 \pm 4.1%, Fig. 1b, c); other markers are shown in Extended Data Fig. 1a. Approximately 30% of the GFP⁺ cells in the CC remained unspecified at this time. None of the GFP-labelled cells expressed neuronal (NeuN) or microglial (Iba1 and CD11b) markers (data not shown). Ten weeks after recovery from cuprizone diet, the numbers of GFP⁺ cells in the demyelinated CC increased from about 17 to 48 cells per section and consisted of PDGFR- α ⁺ OPCs (18.4 \pm 1.8%), CC1⁺ oligodendrocytes (58.3 \pm 5.3%) and GFAP⁺ astrocytes (28.7 \pm 11.6%) (Fig. 1c), accounting for all the GFP⁺ cells in the CC. These results suggest that Gli1⁺ NSCs continue to generate glial cells in the CC for a prolonged period after demyelination. In addition to the cells within the CC, GFP-labelled cells located outside the CC frequently increased with demyelination; these correspond to a subset of protoplasmic astrocytes that are responsive to Shh¹⁷.

The newly generated, NSC-derived oligodendrocytes remyelinated axons, as evidenced by GFP-labelled processes that flanked nodes of Ranvier (Fig. 1e), overlapped with the paranodal marker Caspr (Extended Data Fig. 1b), and co-expressed myelin oligodendrocyte glycoprotein (MOG) and myelin basic protein (MBP), but not the Schwann cell myelin protein P0 (Extended Data Fig. 1b). Immunoelectron microscopy of the CC in *Gli1*^{CreERT2/+} mice crossed to a membrane GFP reporter (*Rosa-TdTomato-mGFP*) demonstrated GFP labelling in compact myelin sheaths surrounding axons (Fig. 1e), corroborating that these NSCs form remyelinating oligodendrocytes.

For comparison, we also fate-mapped the entire pool of NSCs, in healthy brains and after demyelination using a *Nestin*^{CE/+} driver line

¹New York University Neuroscience Institute, Department of Neuroscience and Physiology, New York University School of Medicine, New York, New York 10016, USA. ²The Kimmel Center for Biology and Medicine of the Skirball Institute, New York University School of Medicine, New York, New York 10016, USA.

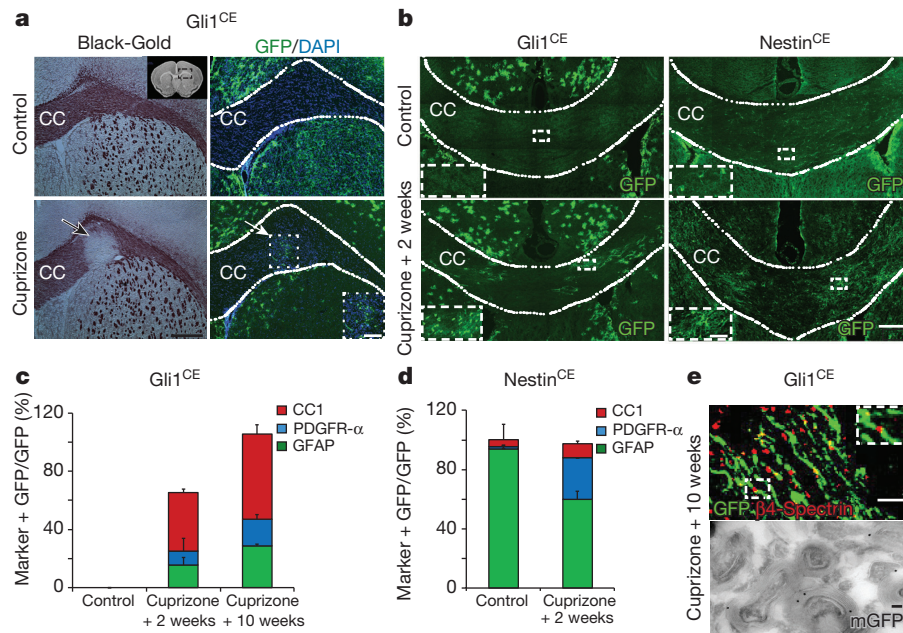


Figure 1 | Gli1-expressing cells are recruited to, and generate myelinating oligodendrocytes at sites of demyelination. **a**, Serial sections from the brains of *Gli1*^{CE/+} mice on a control or cuprizone diet (at the peak of demyelination) were stained with Black-Gold myelin (left) or immunostained for GFP (right). Inset, top left: a coronal section of the forebrain; the box highlights the area of CC analysed. Black arrow, lower left: a demyelinated region of the CC. Labelled cells in the CC are restricted to the site of demyelination (white arrow, lower right). *N* = 5 mice per group. Scale bar, 100 μm; inset, 40 μm. **b**, Comparison of *Gli1*^{CE/+} mice with *Nestin*^{CE/+} mice shows GFP⁺ cells are present in the CC of *Nestin*^{CE/+} mice but not of *Gli1*^{CE/+} mice on control diet (top). During recovery from cuprizone, GFP⁺ cells appear

in the CC of *Gli1*^{CE/+} mice (bottom) and increase in the CC of *Nestin*^{CE/+} mice. *N* = 5 mice per group/genotype. Scale bars, 100 μm; inset, 40 μm.

c, d, Percentage of GFP⁺ oligodendrocyte progenitors (PDGFR-α⁺), mature oligodendrocytes (CC1⁺), and astrocytes (GFAP⁺) in the CC of *Gli1*^{CE/+} (c) and *Nestin*^{CE/+} (d) mice. *N* = 5 mice per group/genotype. **e**, Ten weeks after cessation of cuprizone, GFP-labelled processes flank nodes of Ranvier, demarcated by β4-spectrin (top). *N* = 3 mice. Scale bar, 10 μm.

Immunoelectron microscopy shows membrane-targeted GFP (mGFP), indicated by immunogold particles, present within compact myelin sheaths that surround axons (bottom). *N* = 3 mice. Scale bar, 100 nm. Data are mean ± s.e.m.

that preferentially labels the SVZ¹⁸ (Fig. 1b). We also observed GFP⁺ cells outside the SVZ, including some cells within the healthy CC (Fig. 1b). The cells present in the healthy CC were largely astrocytes (93.81 ± 3.34%); a small proportion were OPCs (1.81 ± 1.57%) and oligodendrocytes (4.42 ± 0.29%) (Fig. 1b, d). Two weeks after recovery from cuprizone, there was a striking increase in the numbers of labelled cells in the CC, consistent with a recent report², associated with an increase in the percentages of OPCs (28.49 ± 18.52%) and oligodendrocytes (9.35 ± 3.25%) and a commensurate reduction in the proportion of astrocytes (59.72 ± 9.54%) (Fig. 1b, d). These results do not distinguish whether the oligodendrocytes present in the CC during remyelination were generated primarily from NSCs, from pre-existing precursors within the CC or both. They emphasize that, in contrast to the broader *Nestin*⁺ pool of NSCs, *Gli1* demarcates a distinct set of NSCs in the SVZ that are recruited only upon demyelination and preferentially fated to oligodendroglia.

To assess whether cells that enter demyelinated lesions in the CC were actively responding to Shh, we analysed *Gli1* expression using *Gli1*^{nLacZ} mice, which express nuclear LacZ from the *Gli1* locus¹⁹. In healthy brains, *Gli1* was expressed by cells in the cortex, basal forebrain and ventral SVZ (Extended Data Fig. 2a, b). NSCs that co-expressed GFAP and *Gli1* were also present in the human SVZ (Extended Data Fig. 2c). Outside the SVZ, labelled cells co-expressed GFAP (data not shown) but not PDGFR-α (Extended Data Fig. 2b) and thus correspond to a subset of mature astrocytes but not OPCs, as previously reported¹⁷. No LacZ⁺ cells were present in the CC during or after recovery from cuprizone- or LPC-mediated demyelination (data not shown), further indicating that Shh-responsive cells do not arise from within the callosum. Thus, although NSCs in the SVZ actively respond to Shh, their progeny, upon entry into the callosum, do not. These results are consistent with the minimal expression of Gli proteins by OPCs or oligodendrocytes during

development (http://web.stanford.edu/group/barres_lab/brain_rnaseq.html) and indicate Shh signalling is decreased in OPCs during both normal development and active remyelination.

A time-course analysis of fate-mapped cells in *Gli1*^{CE/+} mice strongly suggests NSCs emigrate from the SVZ into the CC after demyelination (Extended Data Fig. 2d–f) and argues against trans-differentiation of *Gli1*-expressing astrocytes in the cortex as a source of labelled cells. Together, these data indicate that a Shh-responsive pool of NSCs in the SVZ are recruited specifically to the demyelinated CC where they downregulate *Gli1* and differentiate into mature, myelinating oligodendrocytes.

Downregulation of *Gli1* expression by remyelinating SVZ-derived OPCs raised the possibility that inhibiting Shh signalling might further augment remyelination—a notion that contrasts with previous studies showing that Shh signalling can promote repair after CNS injury²⁰ and is required for remyelination¹⁰. We first analysed the effects of a partial loss of Shh signalling by fate mapping NSCs in *Gli1* null (*Gli1*^{CE/nLacZ};RCE) versus heterozygous (*Gli1*^{CE/+};RCE) mice. Consistent with our hypothesis, there were many more GFP-labelled cells in the remyelinating CC of *Gli1* nulls (63 ± 6.5 per section) than in *Gli1* heterozygotes (17.3 ± 2.6 per section, Fig. 2a, b), and a much higher percentage of these were mature oligodendrocytes in the nulls (81.3 ± 4.4%) than in the heterozygotes (40.2 ± 15.1%, Fig. 2c). Overall, there were ~7.5-fold more GFP-labelled, mature oligodendrocytes in the *Gli1* null versus heterozygotes mice. The proportion of labelled OPCs in the nulls versus heterozygotes was similar (5.22 ± 2.5% versus 9.8 ± 8.7%, Fig. 2c) whereas that of GFAP-expressing astrocytes was significantly reduced (2.6 ± 1.3% versus 15.5 ± 4.1%, Fig. 2c). Finally, overall myelin levels in the CC were increased in the *Gli1* nulls compared with heterozygotes 3 weeks after cessation of cuprizone (Fig. 2d), strongly suggesting this enhanced NSC remyelination is physiologically significant.

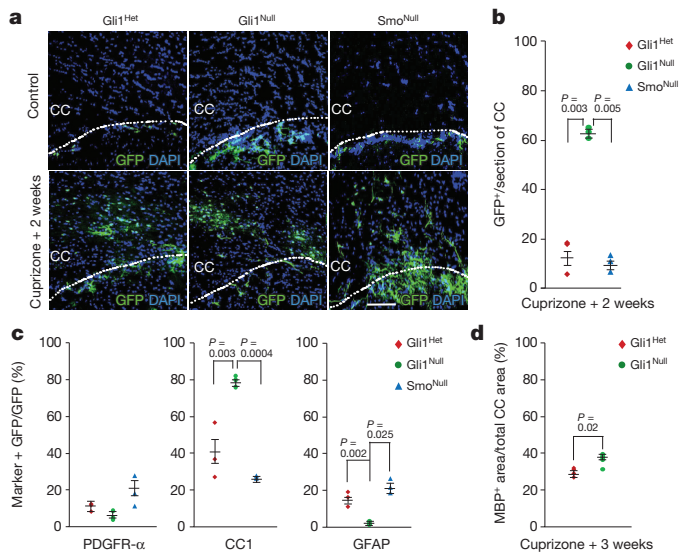


Figure 2 | Loss of Gli1 enhances oligodendroglial differentiation during remyelination. **a**, Brains of *Gli1*^{CE/+} (*Gli1*^{Het}), *Gli1*^{CE/nLacZ} (*Gli1*^{Null}) and *Gli1*^{CE/+}; *Smo*^{fx/fx} (*Smo*^{Null}) mice were analysed 2 weeks after cessation of cuprizone by immunofluorescence. GFP⁺ cells are only observed in the CC of mice receiving cuprizone. *N* = 10 mice per group/genotype. Scale bar, 50 μm. **b**, Quantification of the numbers of GFP⁺ cells in the CC shows a significant increase in *Gli1*^{Null} mice compared with *Gli1*^{Het} and *Smo*^{Null} mice. *N* = 3 mice per group/genotype. **c**, *Gli1*^{Null} mice have a greater proportion of labelled mature oligodendrocytes and reduced proportions of astrocytes than do *Gli1*^{Het} and *Smo*^{Null} mice. *N* = 3 mice per group/genotype. **d**, Three weeks after cessation of cuprizone, *Gli1*^{Null} mice have enhanced MBP expression in the CC. *N* = 4 mice per group/genotype. Scale bar, 50 μm. Data are mean ± s.e.m., Student's *t*-test.

While cuprizone primarily demyelinate the CC, it also has modest effects on other white matter tracts¹⁶. Accordingly, labelled oligodendrocytes in the *Gli1* nulls were also present in white matter tracts at sites distant from the SVZ including the lateral striatum, anterior commissure (Extended Data Fig. 3a, b) and the optic nerve (data not shown). No GFP-labelled cells were seen in the CC or other white matter tracts of heterozygotes or nulls on a control diet (Fig. 2a). Thus, in the adult, this effect of *Gli1* is specific to remyelination.

Myelination also started significantly earlier in the *Gli1* nulls than in the heterozygotes (Extended Data Fig. 4a, b). In addition, the CC were slightly larger on average in the adult *Gli1* nulls versus heterozygotes on the basis of Black-Gold myelin staining (Extended Data Fig. 4c–e), although differences were not statistically significant. Thus, during development, *Gli1* expression delays the onset of myelination and in the adult substantially inhibits remyelination by the Shh-responsive NSCs.

As *Gli1* is expressed only upon sustained, high-level Shh signalling, we asked whether complete abrogation of this pathway also enhanced myelination. In contrast to loss of *Gli1*, loss of canonical Shh signalling by conditional ablation of smoothened in *Gli1*-expressing NSCs (*Gli1*^{CE/+}; *Smo*^{fx/fx}; *RCE* mice) did not increase the numbers of labelled cells in the CC or alter their cell fates (Fig. 2a–c). Thus, loss of *Gli1* is distinct from loss of Shh signalling, emphasizing the specificity of the effects of *Gli1*.

We also examined the effects of activating canonical Shh signalling in the presence or absence of *Gli1*. To this end, we expressed an activated (M2) form of *Smo*²¹ in NSCs upon tamoxifen treatment in *Gli1* heterozygous (*Gli1*^{CE/+}; *Smo*^{M2}) and null (*Gli1*^{CE/nLacZ}; *Smo*^{M2}) mice; the activated *Smo*^{M2} allele is itself fused to yellow fluorescent protein, allowing fate mapping. Again, labelled cells were only detected in the CC in mice that had undergone demyelination (Extended Data Fig. 5a). In the *Gli1*^{CE/+}; *Smo*^{M2} mice, many of the GFP⁺ cells

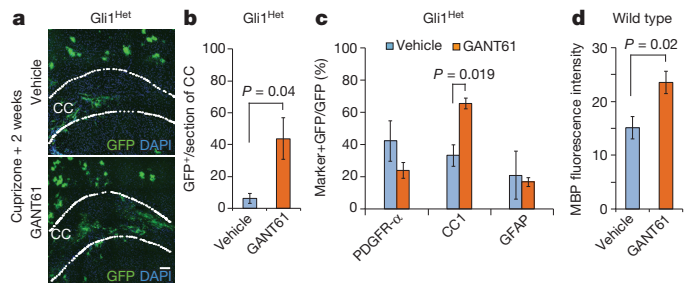


Figure 3 | Pharmacological inhibition of Gli1 promotes NSC recruitment and differentiation during remyelination. **a**, *Gli1*⁺ cells were fate-mapped in cuprizone-fed *Gli1*^{CE/+} (*Gli1*^{Het}) mice treated with vehicle or GANT61 for 4 weeks before analysis. Mice that received GANT61 had more GFP⁺ cells than those receiving vehicle. *N* = 5 mice per group. Scale bar, 50 μm. **b**, **c**, Mice that received GANT61 had an approximately sevenfold increase in the numbers of GFP-labelled cells in the CC (**b**) and a significant increase in the percentage of mature oligodendrocytes (**c**). **d**, Wild-type mice were fed a cuprizone diet and treated with vehicle or GANT61 for 9 weeks. The brains were examined 6 weeks after recovery from cuprizone diet. Mice that received GANT61 showed significantly higher MBP fluorescence intensity levels than vehicle-treated mice. *N* = 5 mice per group. Data are mean ± s.e.m., Student's *t*-test.

(5.6 ± 2.9 per section) in the CC (Extended Data Fig. 5b) were oligodendrocyte progenitors (57.5 ± 23.2%, Extended Data Fig. 5c), a much higher proportion than in *Gli1*^{CE/+} mice (9.8 ± 8.7%). These results agree with previous studies that showed that increasing the levels of Shh in the brain enhances the generation of OPCs but blocks their maturation²². The *Gli1*^{null}; *Smo*^{M2} mice had many more GFP⁺-labelled cells in the CC than the *Gli1*^{het}; *Smo*^{M2} mice (47.2 ± 25.7 versus 5.6 ± 2.9 per section, Extended Data Fig. 5b) and a significantly greater proportion of these were mature oligodendrocytes (52.5 ± 10.5% versus 28.8 ± 3.8%). On average, the *Gli1*^{null}; *Smo*^{M2} mice had ~16 times as many oligodendrocytes derived from the Shh-responsive NSC pool than the *Gli1*^{het}; *Smo*^{M2} mice (Extended Data Fig. 5b, c). These results indicate loss of *Gli1* has an even greater effect in the context of active Shh signalling, promoting robust recruitment and relieving an arrest of NSCs differentiation into oligodendrocytes in the remyelinating CC.

Other analyses revealed that loss of *Gli1* results in a significant increase in proliferation in Shh-responsive SVZ NSCs but only at the onset of demyelination, namely 3 weeks of cuprizone treatment (Extended Data Fig. 6a, b). There was no difference in the proliferation of NSCs of *Gli1* nulls versus heterozygotes on a normal diet, on a cuprizone diet at 4, 5 or 6 weeks or 2 weeks after removal of cuprizone (data not shown). The higher proliferation at 3 weeks may contribute to preservation of the stem-cell pool, which was unchanged in the SVZ of *Gli1* nulls (see Extended Data Fig. 6c). We did not detect any increase in Shh levels with demyelination or a significant difference in Shh levels between *Gli1* heterozygotes versus nulls (Extended Data Fig. 6d–f). Thus, the major effects of loss of *Gli1* in NSCs appear to be enhanced differentiation and recruitment to lesion sites and an increase in proliferation during demyelination.

These findings suggested that *Gli1* might be a useful therapeutic target to promote remyelination. To test this possibility, we infused GANT61, a small molecule inhibitor of Gli1²³ into the lateral ventricle of *Gli1*^{CE/+} mice via a mini-osmotic pump; we corroborated inhibition of *Gli1* by the reduction of its messenger RNA (mRNA) levels by quantitative PCR (qPCR) (Extended Data Fig. 7a). A similar inhibition of *Gli1* mRNA levels was observed when the drug was administered by intraperitoneal injection and oral gavage, indicating GANT61 can cross the blood–brain barrier efficiently. We infused GANT61 in *Gli1*^{CE/+} mice during the last 2 weeks of the cuprizone diet and continued it for an additional 2 weeks off cuprizone for a total of 4 weeks. Mice that received GANT61 versus vehicle had significantly greater numbers of labelled cells (43.8 ± 22.1 versus 6.4 ± 5 per section,

Fig. 3a, b) and a significantly greater proportion of these were oligodendrocytes ($65.4 \pm 7.5\%$ versus $21.4 \pm 18.7\%$, Fig. 3c). Administration of GANT61 was well tolerated and did not deplete the NSCs in the SVZ (Extended Data Fig. 7b, c). Importantly, mice receiving extended treatment with GANT61, namely during the last 3 weeks of the cuprizone diet and 6 weeks thereafter for a total of 9 weeks, had more myelin in the CC than mice similarly treated with vehicle (Fig. 3d). Thus, the Gli1 inhibitor enhanced the recruitment and differentiation of Shh-responsive NSCs into oligodendrocytes at sites of demyelination, promoting remyelination.

These effects of GANT61 are specific to NSCs as there were no effects on remyelination by OPCs on the basis of their fate mapping in *NG2^{CreERT2}* transgenic mice²⁴ crossed to the RCE reporter (Extended Data Fig. 8). The lack of an effect of GANT61 on OPC remyelination is consistent with the absence of Gli1 expression by these cells. These results also indicate that enhanced repair by NSCs resulting from GANT61 treatment does not come at the expense of OPC remyelination: rather, it is additive.

To address the therapeutic potential of inhibiting Gli1, we examined the effects of GANT61 in a relapsing-remitting model of experimental autoimmune encephalitis (RR-EAE), a physiologically relevant model of inflammatory demyelination and remyelination. In this model, RR-EAE is induced by injecting proteolipid protein (PLP) peptide into wild-type SJL mice²⁵; the severity of the clinical phenotype of the initial attack correlates with the extent of spinal cord inflammation, whereas the late-stage neurological disability in later relapses correlates with axonal loss²⁶ probably resulting from cumulative injury including chronic demyelination. GANT61 was administered by daily oral gavage either prophylactically (that is, at the onset of PLP immunization) or therapeutically (that is, at the onset of symptoms). Neither treatment protocol altered the induction of EAE or the severity of the acute attack (about day 12), suggesting GANT61 did not affect the immune response (Fig. 4a). In agreement, Gli1 was not expressed by cells in the lymphocytic or monocytic lineages isolated from the spleen, thymus or liver of healthy (Extended Data Fig. 9) or cuprizone-treated mice (data not shown).

Of note, GANT61 reduced the severity of the first relapse (around day 27) and significantly enhanced functional recovery during and after the second relapse (around day 46) compared with vehicle treatment (Fig. 4a). We therefore examined myelin levels, axon pathology and motor neuron numbers in the (prophylactic, therapeutic) GANT61- and vehicle-treated lumbar spinal cords at the end of the second relapse phase (that is, day 53). In electron microscopy images, all three groups had significant spinal cord pathology, with disruption of fascicles of myelinated axons and pathology of individual myelinated axons most evident in vehicle-treated compared with the GANT61-treated groups (Fig. 4b, c). In all EAE groups, pathology was most pronounced in small-diameter axons ($<0.5 \mu\text{m}$, Extended Data Fig. 10a–c) which were increased in numbers and had significantly lower G ratios (ratio of axon diameter to myelinated axon diameter); the reduction in axon diameter and corresponding increase in G ratios are probably due to axonal atrophy and suggestive of active demyelination²⁷. We did not detect significant numbers of unmyelinated or thinly myelinated axons in any group (Extended Data Fig. 10a–c). While thin myelin sheaths have long been considered the hallmark of remyelination¹, recent studies suggest remyelination performed by neural stem cells² or in the spinal cord can be of normal thickness²⁸. Analysis of MBP levels supported significantly higher levels of myelin in the drug- versus vehicle-treated groups (Fig. 4f).

The numbers of lower motor neurons, identified by co-staining for NeuN and choline acetyltransferase (ChAT)²⁹, were significantly reduced in the spinal cords of vehicle-treated (16.93 ± 1.92 per section) versus healthy control mice (27.59 ± 5.9 per section, Fig. 4d, e). Numbers of lower motor neurons in the GANT61-treated groups were intermediate and closer to controls (prophylactic: 23.4 ± 4.3 per section; therapeutic: 26.93 ± 8.4 per section), providing strong presumptive evidence of neural protection, potentially because of reduced axon pathology and/or remyelination. The relative preservation of lower motor neurons is consistent with the reduction in axonal pathology and probably accounts for improved functional outcomes. Taken together, these results suggest inhibition of Gli1 enhances remyelination and thereby protects neurons from degeneration in this EAE

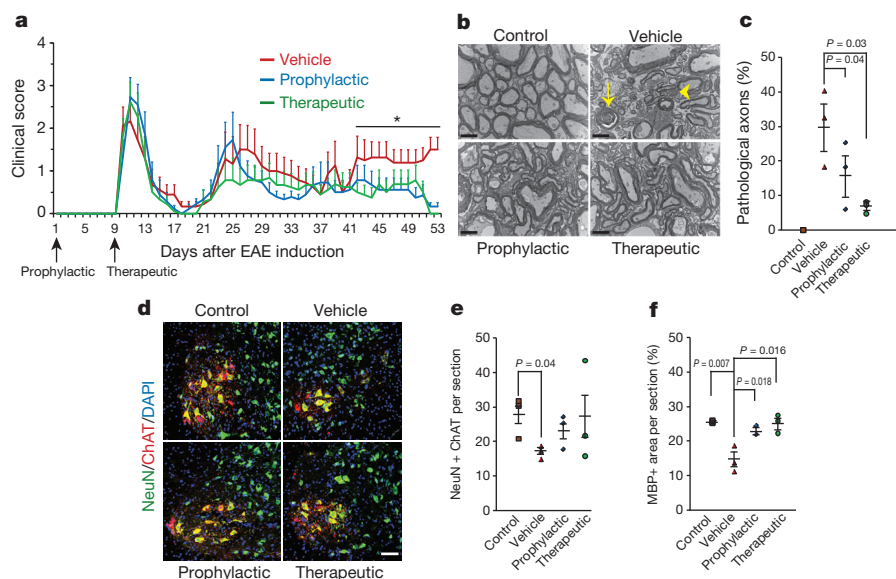


Figure 4 | GANT61 improves functional outcomes and is neuroprotective in a RR-EAE model. **a**, EAE clinical scores after prophylactic or therapeutic treatment with GANT61 compared with vehicle administration. $N = 9$ mice per group. Data are mean \pm s.e.m., * $P < 0.05$, two-way analysis of variance with Tukey's multiple comparison test. **b**, **c**, Electron micrographs from the ventral lumbar spinal cords of control, GANT61- and vehicle-treated EAE mice shows axonal pathology (**b**) including axolysis (arrow) and dense axoplasm (arrowhead). Quantification of 500 axons per group from three mice per

group (**c**) shows a higher proportion of pathological axons in vehicle-treated group compared with GANT61-treated EAE mice. Scale bar, 2 nm. $N = 3$ mice per group. **d**, **e**, Alpha motor neurons labelled for NeuN and ChAT in the lumbar spinal cords of control, vehicle- and GANT61-treated EAE mice (**d**). Quantification shows more motor neurons in the GANT61- versus vehicle-treated mice (**e**). Scale bar, 50 μm . $N = 3$ mice per group. **f**, Quantification for the area of MBP expression shows lowest levels of myelin in the vehicle-treated group. $N = 3$ mice per group. Data are mean \pm s.e.m., Student's *t*-test.

model without altering the immune response. Gli1 is primarily expressed in the spinal cord by a subset of astrocytes within the parenchymal grey matter and by cells located around the central canal, a site of spinal cord NSCs³⁰ but not by NG2⁺ OPCs (Extended Data Fig. 10d). These findings suggest the effects of inhibiting Gli1 on myelination and neuroprotection are mediated by direct effects on neural stem cells or potentially by indirect effects mediated by parenchymal astrocytes.

These studies demonstrate that inhibiting Gli1 appears to be a well-tolerated and effective strategy for mobilizing and enhancing the differentiation of a resident population of Shh-responsive neural stem cells. The findings of elevated levels of myelin, reduced axon pathology and preservation of lower motor neurons in a RR-EAE model further support its therapeutic potential. This approach may therefore be useful in aiding repair in multiple sclerosis and other demyelinating neurological disorders. The findings also highlight that different pools of remyelinating cells (that is, NSCs and OPCs) might require different strategies to promote their remyelination of axons, raising the possibility of using combinatorial strategies to enhance repair.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 11 April; accepted 15 July 2015.

Published online 30 September 2015.

- Franklin, R. J. M. & Goldman, S. A. Glia disease and repair—remyelination. *Cold Spring Harb. Perspect. Biol.* <http://dx.doi.org/10.1101/cshperspect.a020594> (2015).
- Xing, Y. L. *et al.* Adult neural precursor cells from the subventricular zone contribute significantly to oligodendrocyte regeneration and remyelination. *J. Neurosci.* **34**, 14128–14146 (2014).
- Scolding, N. *et al.* Oligodendrocyte progenitors are present in the normal adult human CNS and in the lesions of multiple sclerosis. *Brain* **121**, 2221–2228 (1998).
- Gensert, J. M. & Goldman, J. E. Endogenous progenitors remyelinate demyelinated axons in the adult CNS. *Neuron* **19**, 197–203 (1997).
- Zawadzka, M. *et al.* CNS-resident glial progenitor/stem cells produce Schwann cells as well as oligodendrocytes during repair of CNS demyelination. *Cell Stem Cell* **6**, 578–590 (2010).
- Menn, B. *et al.* Origin of oligodendrocytes in the subventricular zone of the adult brain. *J. Neurosci.* **26**, 7907–7918 (2006).
- Nait-Oumesmar, B. *et al.* Activation of the subventricular zone in multiple sclerosis: evidence for early glial progenitors. *Proc. Natl Acad. Sci. USA* **104**, 4694–4699 (2007).
- Fuccillo, M., Joyner, A. L. & Fishell, G. Morphogen to mitogen: the multiple roles of hedgehog signalling in vertebrate neural development. *Nature Rev. Neurosci.* **7**, 772–783 (2006).
- Petrova, R. & Joyner, A. L. Roles for Hedgehog signaling in adult organ homeostasis and repair. *Development* **141**, 3445–3457 (2014).
- Ferent, J., Zimmer, C., Durbec, P., Ruat, M. & Traiffort, E. Sonic Hedgehog signaling is a positive oligodendrocyte regulator during demyelination. *J. Neurosci.* **33**, 1759–1772 (2013).
- Ingham, P. W. & McMahon, A. P. Hedgehog signaling in animal development: paradigms and principles. *Genes Dev.* **15**, 3059–3087 (2001).
- Ahn, S. & Joyner, A. L. Dynamic changes in the response of cells to positive hedgehog signaling during mouse limb patterning. *Cell* **118**, 505–516 (2004).
- Dessaud, E., McMahon, A. P. & Briscoe, J. Pattern formation in the vertebrate neural tube: a sonic hedgehog morphogen-regulated transcriptional network. *Development* **135**, 2489–2503 (2008).
- Ahn, S. & Joyner, A. L. *In vivo* analysis of quiescent adult neural stem cells responding to Sonic hedgehog. *Nature* **437**, 894–897 (2005).
- Sousa, V. H., Miyoshi, G., Hjerling-Leffler, J., Karayannis, T. & Fishell, G. Characterization of Nkx6-2-derived neocortical interneuron lineages. *Cereb. Cortex* **19** (Suppl. 1), i1–i10 (2009).
- Matsushima, G. K. & Morell, P. The neurotoxicant, cuprizone, as a model to study demyelination and remyelination in the central nervous system. *Brain Pathol.* **11**, 107–116 (2001).
- Garcia, A. D., Petrova, R., Eng, L. & Joyner, A. L. Sonic hedgehog regulates discrete populations of astrocytes in the adult mouse forebrain. *J. Neurosci.* **30**, 13597–13608 (2010).
- Balordi, F. & Fishell, G. Mosaic removal of hedgehog signaling in the adult SVZ reveals that the residual wild-type stem cells have a limited capacity for self-renewal. *J. Neurosci.* **27**, 14248–14259 (2007).
- Bai, C. B., Auerbach, W., Lee, J. S., Stephen, D. & Joyner, A. L. Gli2, but not Gli1, is required for initial Shh signaling and ectopic activation of the Shh pathway. *Development* **129**, 4753–4761 (2002).
- Bambakidis, N. C. & Onwuzulike, K. Sonic Hedgehog signaling and potential therapeutic indications. *Vitam. Horm.* **88**, 379–394 (2012).
- Xie, J. *et al.* Activating *Smoothed* mutations in sporadic basal-cell carcinoma. *Nature* **391**, 90–92 (1998).
- Rowitch, D. H. Sonic hedgehog regulates proliferation and inhibits differentiation of CNS precursor cells. *J. Neurosci.* **19**, 8954–8965 (1999).
- Lauth, M., Bergstrom, A., Shimokawa, T. & Toftgard, R. Inhibition of GLI-mediated transcription and tumor cell growth by small-molecule antagonists. *Proc. Natl Acad. Sci. USA* **104**, 8455–8460 (2007).
- Zhu, X. *et al.* Age-dependent fate and lineage restriction of single NG2 cells. *Development* **138**, 745–753 (2011).
- Tuohy, V. K., Sobel, R. A. & Lees, M. B. Myelin proteolipid protein-induced experimental allergic encephalomyelitis. Variations of disease expression in different strains of mice. *J. Immunol.* **140**, 1868–1873 (1988).
- Wujek, J. R. *et al.* Axon loss in the spinal cord determines permanent neurological disability in an animal model of multiple sclerosis. *J. Neuropathol. Exp. Neurol.* **61**, 23–32 (2002).
- Recks, M. S. *et al.* Early axonal damage and progressive myelin pathology define the kinetics of CNS histopathology in a mouse model of multiple sclerosis. *Clin. Immunol.* **149**, 32–45 (2013).
- Powers, B. E. *et al.* Remyelination reporter reveals prolonged refinement of spontaneously regenerated myelin. *Proc. Natl Acad. Sci. USA* **110**, 4075–4080 (2013).
- Aharoni, R. *et al.* Distinct pathological patterns in relapsing-remitting and chronic models of experimental autoimmune encephalomyelitis and the neuroprotective effect of glatiramer acetate. *J. Autoimmun.* **37**, 228–241 (2011).
- Barnabe-Heider, F. *et al.* Origin of new glial cells in intact and injured adult spinal cord. *Cell Stem Cell* **7**, 470–482 (2010).

Supplementary Information is available in the online version of the paper.

Acknowledgements We thank A. Joyner for providing mouse lines and for advice during the course of this project, M. Bhat, M. Wegner, and M. Rasband for providing antibodies, G. Zanazzi for providing human brain tissue, A. Liang for assistance with immunoelectron microscopy, and G. Multani for technical assistance during initial studies. This research was supported by grants to J.L.S. from the New York State Department of Health Stem Cell Board and the National Multiple Sclerosis Society. J.S. was a recipient of a postdoctoral fellowship from the National Multiple Sclerosis Society.

Author Contributions J.S. performed the experiments, analysed the data and co-wrote the paper with J.L.S. All authors contributed to the design of individual experiments, reviewed individual results and assisted with portions of manuscript preparation.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare competing financial interests: details are available in the online version of the paper. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to J.S. (jayshree.samanta@nyumc.org) or J.L.S. (james.salzer@nyumc.org).

METHODS

No statistical methods were used to predetermine sample size.

Fate mapping and demyelination. Ten-week-old mice were administered 5 mg tamoxifen (Sigma) in corn oil on alternate days for a total of four intraperitoneal injections. *Nestin*^{CE/+}, *Gli1*^{CE/+}, *Gli1*^{CE/nLacZ}, *Gli1*^{nLacZ/+}, *Gli1*^{nLacZ/nLacZ} and *Gli1*^{CE/+}; *Smo*^{fx/fx} mice, maintained on Swiss-Webster background, were fed 0.4% cuprizone³¹. *Gli1*^{CE/+}; *Smo*^{M2}, *Gli1*^{CE/nLacZ}; *Smo*^{M2}, *NG2*^{CE} and *Shh*^{CE/+} mice were maintained in C57Bl/6 background and fed 0.2% cuprizone diet to obtain comparable demyelination in the CC³¹ (Supplementary Table 1). No labelling was seen in the absence of tamoxifen administration. For the LPC model, 2 µl of 1% LPC (Calbiochem) were stereotactically injected into the CC at 1.5 mm anterior, 1.2 mm lateral and 2.2 mm ventral to the bregma.

Drug administration. GANT61 (Enzo Life Sciences) (5.6 mg kg⁻¹ per day) or vehicle (50% DMSO + 15% ethanol in PEG400) were delivered via mini-osmotic pump (model 2004, Durect) at a rate of 0.25 µl h⁻¹ for 4 weeks into the lateral ventricle of *Gli1*^{CE/+} mice at 0.5 mm anterior, 0.75 mm lateral and 2.5 mm ventral to the bregma. For intraperitoneal and oral routes, 50 mg kg⁻¹ GANT61 dissolved in ethanol:corn oil (1:4) was administered daily for 4 weeks.

Immunostaining. For Black-Gold (Millipore) myelin staining, mice were perfused with 4% PFA, and 20 µm coronal cryosections were stained according to the manufacturer's protocol. For all other analyses, mice were perfused with Prefer (Anatech) and 20 µm coronal cryosections were processed for immunofluorescence with rabbit or chicken anti-GFP (1:1,000, Invitrogen) and one of the following antibodies: rat anti-PDGFR-α (1:200, BD Biosciences); rabbit anti-NG2 (1:200, Millipore), anti-S100β (1:600, Dako) and β4-spectrin (1:4,000, from M. Rasband); mouse anti-CC1 (1:400, Calbiochem), anti-GFAP (1:400, Sigma), anti-NeuN (1:200, Millipore), anti-Shh (1:500, DSHB), anti-MOG (1:50, Sigma) and anti-MBP (1:500, Millipore); goat anti-LacZ (1:2,000, Biolegend) and anti-ChAT (1:200, Millipore); chicken anti-P0 (1:200, Millipore); guinea pig anti-Caspr (1:3,000, from M. Bhat) and anti-Sox10 (1:1,000, from M. Wegner). Cryosections (30 µm) of an autopsy specimen of a healthy human brain (provided by New York Brain Bank of Columbia University) were stained with rabbit anti-Gli1 (1:2,000, Abcam) and mouse anti-GFAP (1:400, Sigma). Secondary antibodies were goat or donkey anti-species conjugated with Alexafluors (1:1,000, Molecular Probes). Nuclei were counterstained with Hoechst 33258 (1:5,000, Invitrogen). Fluorescent images were obtained as Z-stacks of 1 µm optical sections using a confocal laser-scanning microscope (LSM 510, Zeiss) and processed using Adobe Photoshop. At least ten sections per mouse were analysed and data from three to five mice were combined to determine the average and standard deviation. The investigators were blinded to allocation during experiments and outcome assessment. Student's *t*-test was performed to calculate *P* values.

Flow cytometry. Cells from the thymus, spleen and liver were stained with the following conjugated antibodies: CD4 (RM4-5), CD8 (53-6.7), CD19 (6D5), Ly6c (HK1.4), CD11b (M1/70), CD11c (N418), CD45 (30-F11), PDCA-1 (927), B220 (RA3-6B2) and CD3 (145-2C11). Antibodies were purchased from eBioscience or Biolegend and the staining was performed according to the manufacturer's instructions. DAPI (Invitrogen) was always used to exclude dead cells from the analysis. Stained cells were analysed on a LSRII flow cytometer (BD) and data processed using FlowJo (Tree Star). The investigators were blinded to the genotype of the cells being sorted.

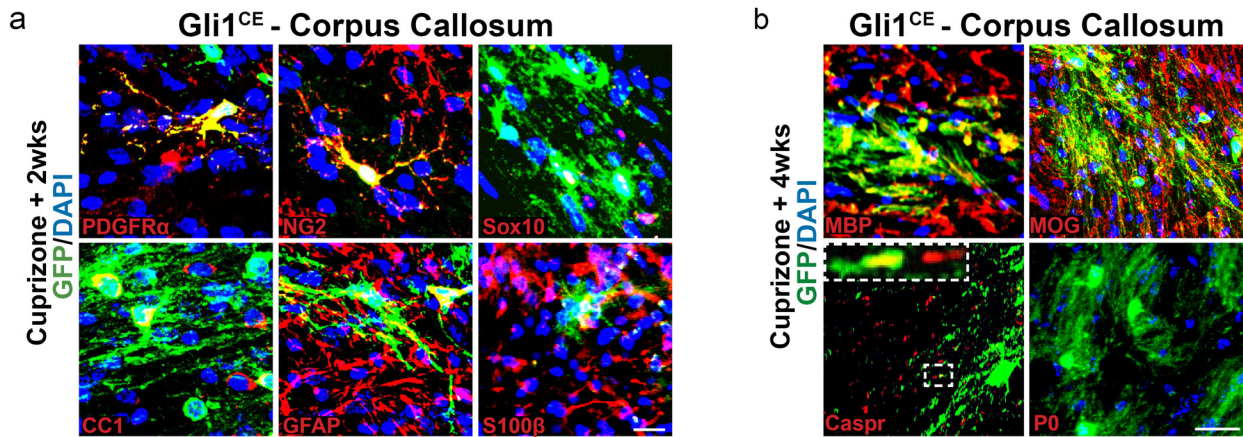
EdU labelling. EdU (Invitrogen; 200 mg kg⁻¹) dissolved in saline was administered by intraperitoneal injection 1 day before harvesting the brains. EdU was visualized using the AlexaFluor-594 Click-iT EdU Cell Proliferation Assay Kit (Invitrogen) according to the manufacturer's instructions.

EAE model. The EAE experiments were randomized. The PLP-induced EAE model was performed and scored using the EAE induction kit (Hooke Labs) on 8-week-old female SJL mice according to the manufacturer's instructions. We used a random number generator to assign each mouse to a group and used an inclusion criteria of clinical score of 2 or above for final analysis (*N* = 9 per group). Mice received either vehicle (ethanol:corn oil (1:4)) or 50 mg kg⁻¹ GANT61 via oral gavage until the end of the experiment, namely day 53 after induction. Separate investigators performed the clinical scoring and the drug treatments and were blinded to each other; the histological analysis was also done blinded.

Electron microscopy. Immunoelectron microscopy was performed on brain sections to detect mGFP in the myelin wraps around axons in the CC according to a previous study³². The investigators were blinded to allocation during experiments and outcome assessment.

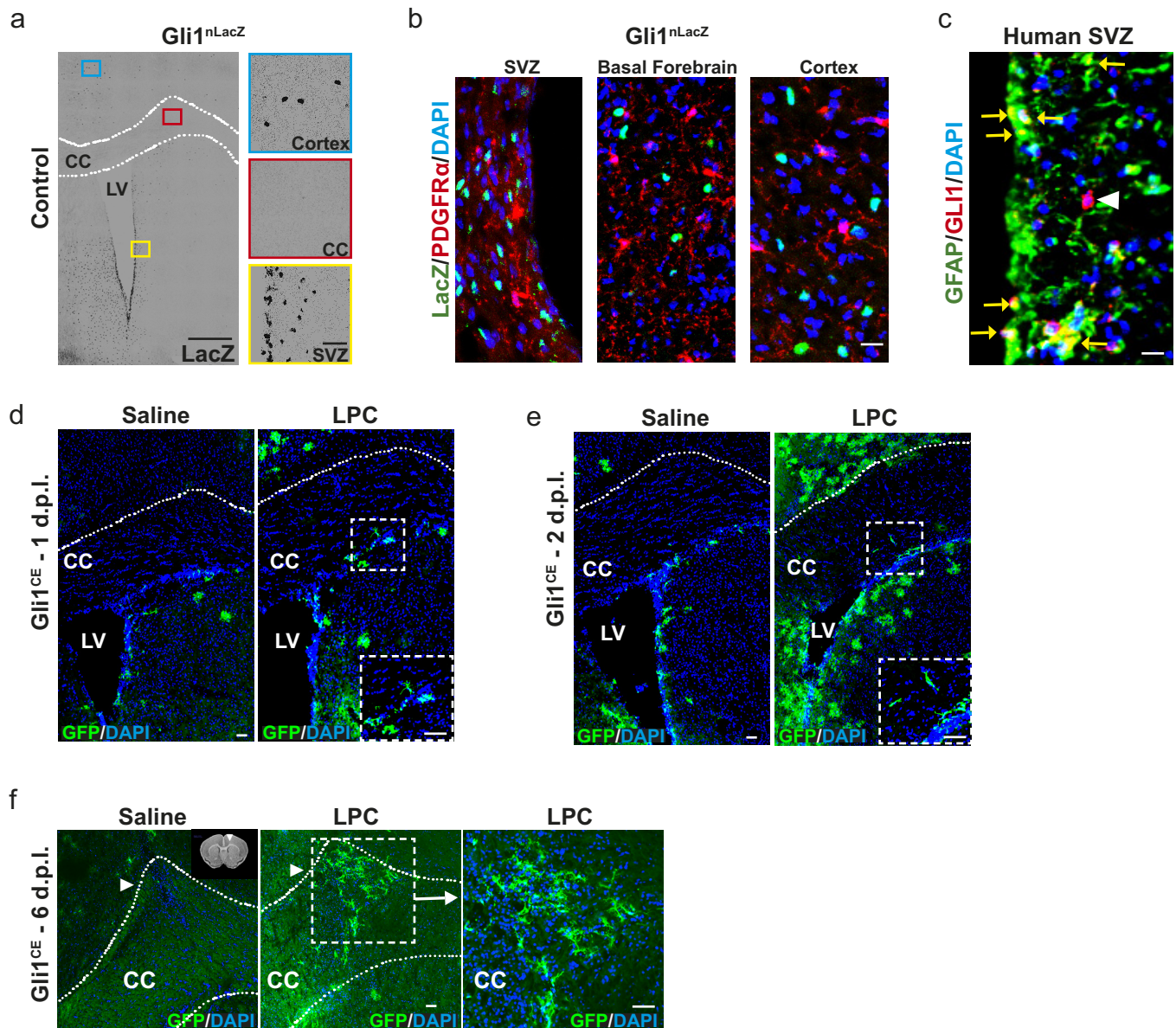
qPCR. mRNA was extracted from the forebrains of three mice in each group using RNeasy kit (Qiagen) and reverse-transcribed to complementary DNA using quantitect kit (Qiagen). Power SYBR Green QPCR Master Mix (Applied Biosystems) was used to perform qPCRs in a Stratagene MX3000P thermal cycler. Primers used were *Gli1* (forward, 5'-CCC ATA GGG TCT CGG GGT CTC AAA C-3'; reverse, 5'-GGA GGA CCT GCG GCT GAC TGT GTA A-3'), *Gli2* (forward, 5'-AGA GAC AGC AGA AGC TAT GCC CAA-3'; reverse, 5'-TGG GCA GCC TCC ATT CTG TTC ATA-3') and *GAPDH* (forward, 5'-GGT GTG AAC GGA TTT GGC CGT ATT G-3'; reverse, 5'-CCG TTG AAT TTG CCG TGA GTG GAG T-3'). The 2^{-ΔΔC_T} method was used to analyse the relative gene expressions and Student's *t*-test was used to calculate *P* values.

31. Elsworth, S. & Howell, J. M. Variation in the response of mice to cuprizone. *Res. Vet. Sci.* **14**, 385–387 (1973).
32. Glausier, J. R., Khan, Z. U. & Muly, E. C. Dopamine D1 and D5 receptors are localized to discrete populations of interneurons in primate prefrontal cortex. *Cereb. Cortex* **19**, 1820–1834 (2009).



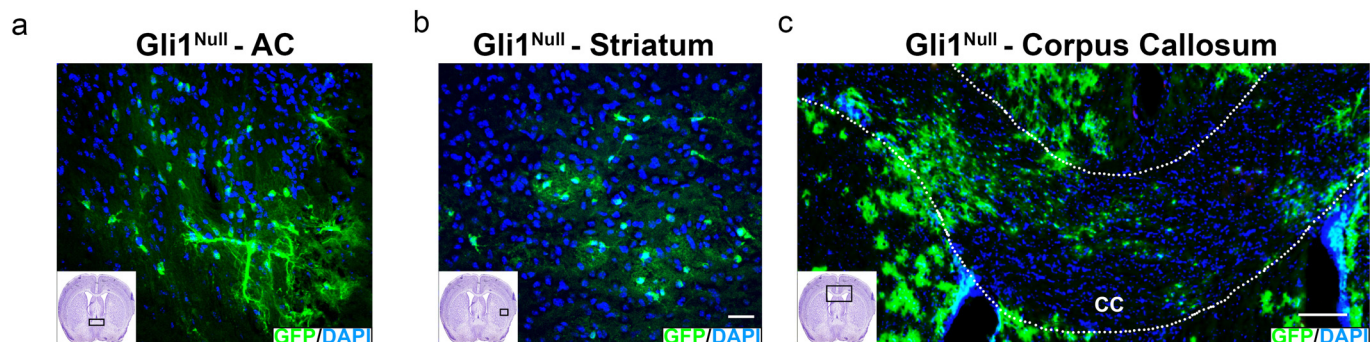
Extended Data Figure 1 | Gli1-expressing cells generate oligodendrocytes after demyelination. Additional markers used to analyse fate-mapped Shh-responsive NSCs are shown. **a**, Two weeks after removal from cuprizone diet, GFP-labelled cells in the CC of Gli1^{CE/+} mice co-expressed the oligodendrocyte progenitor markers PDGFR- α , NG2 and Sox10, the mature

oligodendrocyte marker CC1, and the astrocytic markers GFAP and S100 β . Scale bar, 10 μ m. **b**, Four weeks after removal from cuprizone diet, GFP-labelled processes co-localized with myelin proteins MBP and MOG but not with peripheral myelin protein P0. GFP-labelled processes also overlaid the axonal, paranodal marker Caspr. *N* = 5 mice per group. Scale bar, 10 μ m.



Extended Data Figure 2 | Gli1-expressing neural stem cells in the SVZ egress and generate labelled cells in the CC. **a**, Expression of Gli1 in the forebrain was confirmed in *Gli1^{nLacZ/+}* mice by immunofluorescence for LacZ. Labelled cells were observed in the SVZ, cortex and basal forebrain but not in the CC of mice. The right panel shows the magnified images for the corresponding boxes in the left panel. *N* = 5 mice. Scale bar, 50 μ m. **b**, Double staining for PDGFR- α and LacZ in the *Gli1^{nLacZ/+}* forebrain does not show any co-labelled cells, indicating that Gli1 is not expressed by OPCs. *N* = 5 mice. Scale bar, 50 μ m. **c**, The ventral SVZ lining the body of the lateral ventricle from a human brain specimen shows co-localization of Gli1 with GFAP⁺ cells

(yellow arrow) as well as a Gli1⁺ cell not expressing GFAP (arrowhead). *N* = 1 brain. Scale bar, 50 μ m. **d–f**, Time-course analysis of the SVZ and CC of *Gli1^{CE/+}* mice stereotactically injected with saline (control, left panels) or LPC (right panels) to induce demyelination. *N* = 3 mice per group. No labelled cells were seen within the CC after saline injections; areas of ingress into the LPC-injected CC are boxed. At 1 day post-lesion (d.p.i.) (**d**), GFP-labelled cells diverted towards the CC; at 2 d.p.i. (**e**), a few labelled cells were seen within the CC; at 6 d.p.i. (**f**), many GFP⁺ cells had accumulated at the site of LPC injection (arrowhead). Scale bar, 50 μ m.

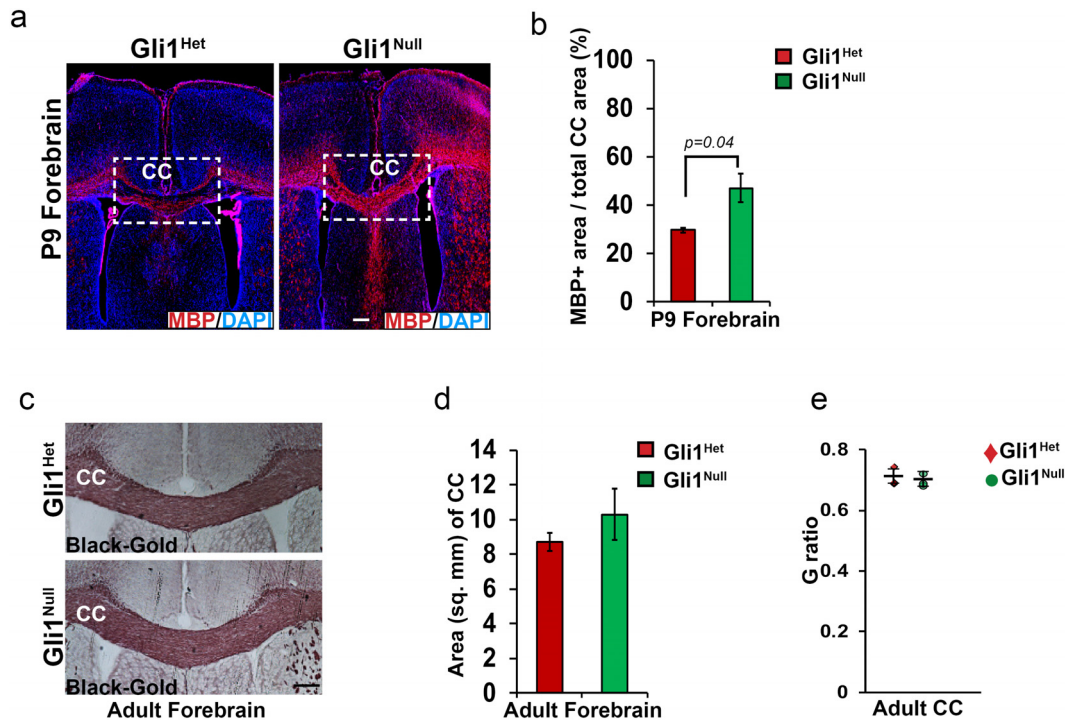


Extended Data Figure 3 | Neural stem cells generate oligodendrocytes in various white matter tracts in Gli1-null brains upon demyelination.

a, b, Fate-mapped Gli1-null cells migrated into the anterior commissure (AC) (**a**) and striatum (**b**) after cuprizone-mediated demyelination (inset shows

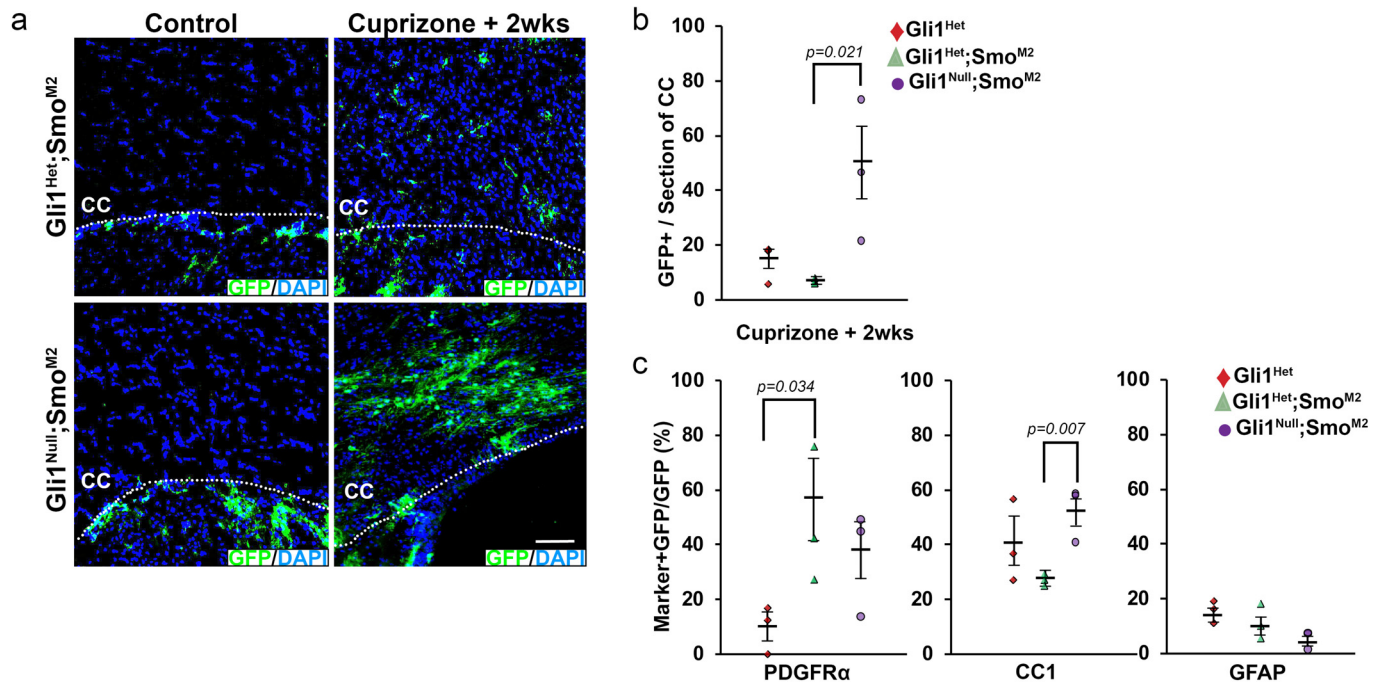
the area of the forebrain) in Gli1^{Null} (*Gli1^{CE/nLacZ}*) mice. Scale bar, 50 μ m.

c, Two weeks after removal from cuprizone diet, GFP-labelled cells are present throughout the CC of Gli1^{Null} (*Gli1^{CE/nLacZ}*) mice. *N* = 5 mice per group. Scale bar, 100 μ m.



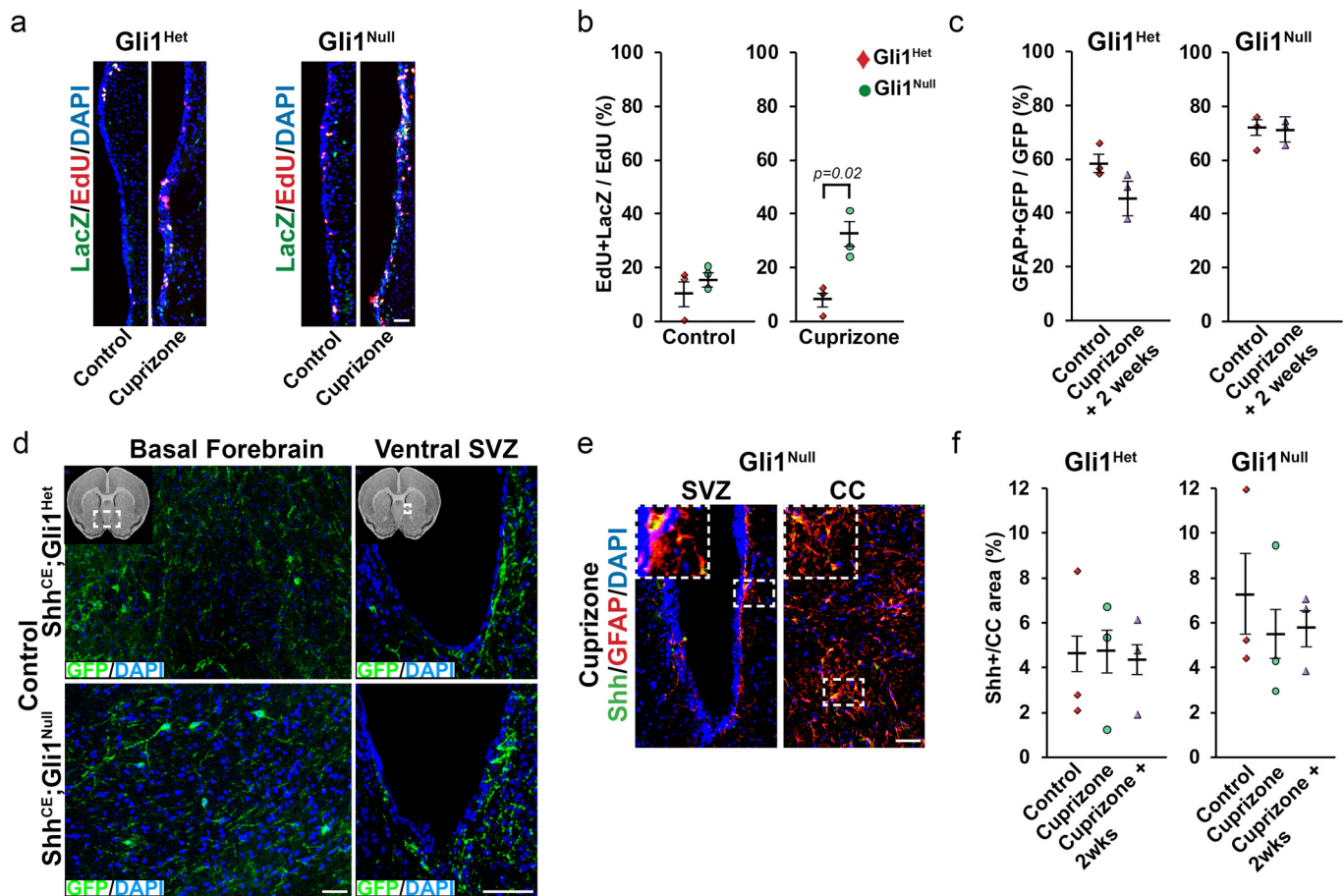
Extended Data Figure 4 | Myelination starts earlier in developing Gli1 null mice. **a**, *Gli1^{nLacZ/nLacZ}* (*Gli1^{Null}*) mice (right panel) show increased MBP levels in the forebrain at postnatal day 9 (P9) compared with *Gli1^{nLacZ/+}* (*Gli1^{Het}*) mice. **b**, Quantification of the extent of MBP expression at P9 in the CC of Gli1 nulls ($47.12 \pm 10.44\%$) versus heterozygotes ($29.71 \pm 1.77\%$) corroborates that myelination is accelerated. $N = 5$ mice per genotype. **c**, Analysis of healthy adult forebrain shows the intensity of Black-Gold myelin stain in Gli1 heterozygotes and nulls was comparable. **d**, Quantification of the

sizes of the CC shows the CC in *Gli1^{Null}* (*Gli1^{nLacZ/nLacZ}*) was slightly larger on average than that in *Gli1^{Het}* (*Gli1^{nLacZ/+}*) mice, although the difference was not statistically significant. $N = 5$ mice per genotype. **e**, Quantification of G ratios from electron micrographs of healthy *Gli1^{Het}* (*Gli1^{nLacZ/+}*) and *Gli1^{Null}* (*Gli1^{nLacZ/nLacZ}*) mice revealed no difference in the thickness of myelin sheaths in the CC. $N = 3$ mice per genotype. Scale bar, 50 μm . Data are mean \pm s.e.m., Student's *t*-test.



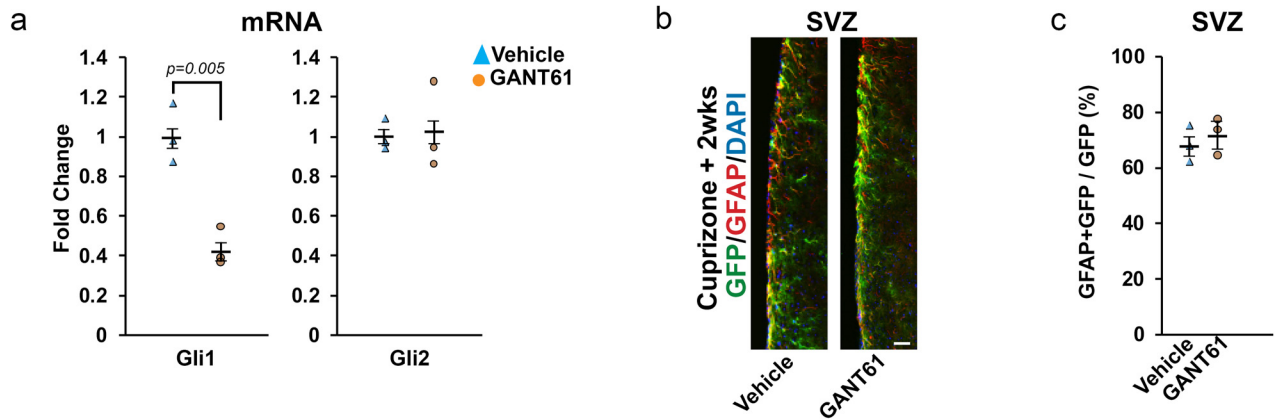
Extended Data Figure 5 | Effects of gain of smoothened function in Gli1-het versus Gli1-null cells during remyelination. Gli1^{Het} (*Gli1*^{CE/+}), Gli1^{Het};Smo^{M2} (*Gli1*^{CE/+};Smo^{M2}) and Gli1^{Null};Smo^{M2} (*Gli1*^{CE/nLacZ};Smo^{M2}) mice were injected with tamoxifen, fed either a regular or a cuprizone-supplemented diet for 6 weeks and analysed by immunofluorescence 2 weeks after removal of cuprizone. **a**, GFP⁺ cells are only seen in the CC of mice on cuprizone diet (right) and not in the control mice (left). **b**, Quantification of the GFP⁺ cells in the CC shows significantly higher numbers of cells in

Gli1^{Null};Smo^{M2} mice compared with Gli1^{Het} and Gli1^{Het};Smo^{M2} mice. **c**, Quantification of the proportion of GFP-labelled co-expressing glial markers in the CC of cuprizone-treated Gli1^{Het}, Gli1^{Het};Smo^{M2} and Gli1^{Null};Smo^{M2} mice shows an increase in percentage of GFP-labelled OPCs (PDGFR- α ⁺) in Gli1^{Het};Smo^{M2} mice and mature oligodendrocytes (CC1⁺) in Gli1^{Null};Smo^{M2} mice. *N* = 3 mice per group for each genotype. Data are mean \pm s.e.m., Student's *t*-test.



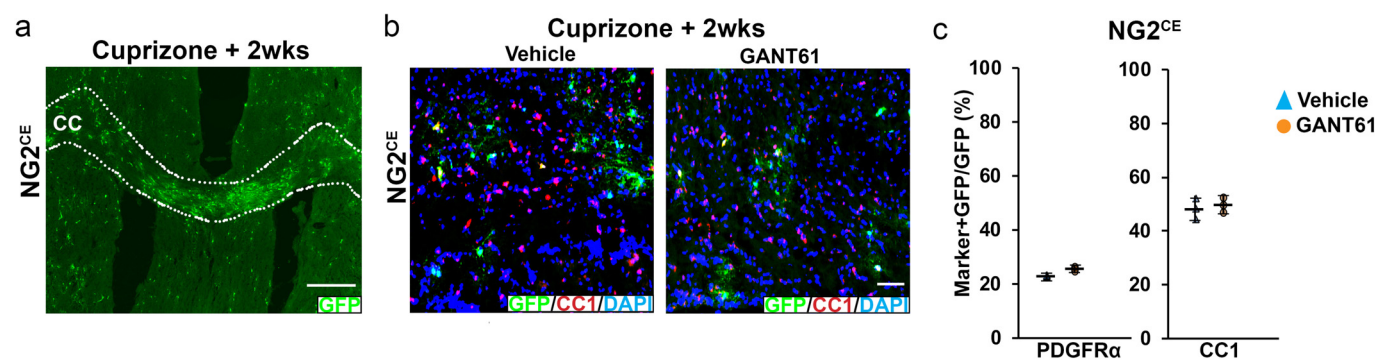
Extended Data Figure 6 | Proliferation of NSCs and expression of Shh in Gli1-null mice. **a, b,** At the start of demyelination (3 weeks of cuprizone diet), Gli1^{Null} (*Gli1^{nLacZ/nLacZ}*) brains have a higher proportion of proliferating nLacZ⁺ neural stem cells indicated by the percentage of EdU-incorporating cells co-expressing nLacZ in the SVZ compared with Gli1^{Het} (*Gli1^{nLacZ/+}*) brains, namely $31 \pm 8.9\%$ versus $8.22 \pm 5.33\%$, respectively. $N = 3$ mice per group for each genotype. Scale bar, 50 μm . **c,** The numbers of fate-mapped Gli1⁺ neural stem cells in the SVZ were quantified as the proportion of GFP⁺ cells co-expressing GFAP in Gli1^{Het} (*Gli1^{CE/+}*) and Gli1^{Null} (*Gli1^{CE/nLacZ}*) mice at 2 weeks of recovery from cuprizone diet. The percentage of GFAP⁺GFP⁺ cells in the SVZ of mice receiving cuprizone diet was comparable to those on a control diet, suggesting that the stem cell pool is not depleted during

remyelination. $N = 3$ mice per group for each genotype. Data are mean \pm s.e.m., Student's *t*-test. **d,** Fate-mapping of Shh expressing cells using an mGFP reporter labels neurons in the basal forebrain (left) with their neurites reaching the ventral SVZ (right) in Gli1^{Het} and Gli1^{Null} brains. **e,** Immunostaining of Gli1^{Null} mice shows Shh in the SVZ and CC after demyelination is mostly co-localized to GFAP-expressing cells. Thus, Shh is produced by neurons of the basal forebrain and binds to a responsive set of astrocytes and neural stem cells. **f,** Quantification of the proportional area of the CC expressing Shh does not show any significant difference between Gli1^{Het} and Gli1^{Null} mice either on control or cuprizone diet. $N = 3$ mice per group for each genotype. Data are mean \pm s.e.m., Student's *t*-test. Scale bar, 50 μm .



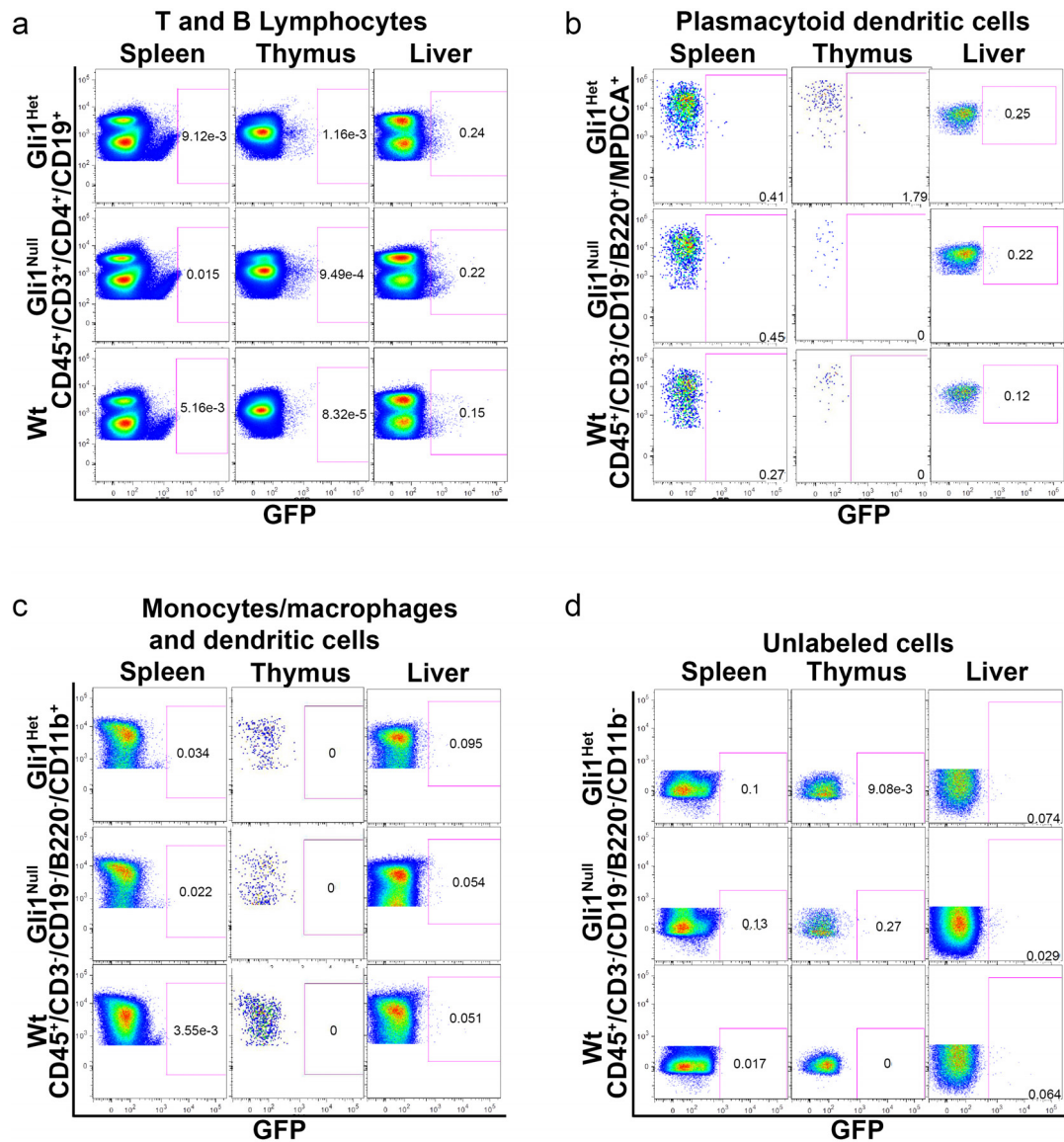
Extended Data Figure 7 | GANT61 reduces Gli1 levels but does not deplete neural stem cells in the SVZ. **a**, Relative expression of *Gli1* and *Gli2* mRNA in the forebrain of *Gli1^{CE/+}* mice was examined by qPCR after administration of GANT61 (50 mg/kg/day) for 4 weeks. GANT61 decreases the mRNA levels of *Gli1* significantly without changing *Gli2* levels. $N = 3$ mice per group. Data are mean \pm s.e.m., Student's *t*-test. **b**, **c**, The numbers of fate-mapped *Gli1*⁺ neural stem cells in the SVZ were analysed by immunofluorescence as the

proportion of GFP⁺ cells co-expressing GFAP in *Gli1^{CE/+}* mice treated with vehicle or GANT61 at 2 weeks of recovery from cuprizone diet (**b**). The percentage of GFAP⁺GFP⁺ cells in the SVZ of mice treated with GANT61 was comparable to those treated with vehicle, suggesting that the stem cell pool is not depleted by GANT61 (**c**). $N = 3$ mice per group. Data are mean \pm s.e.m., Student's *t*-test.



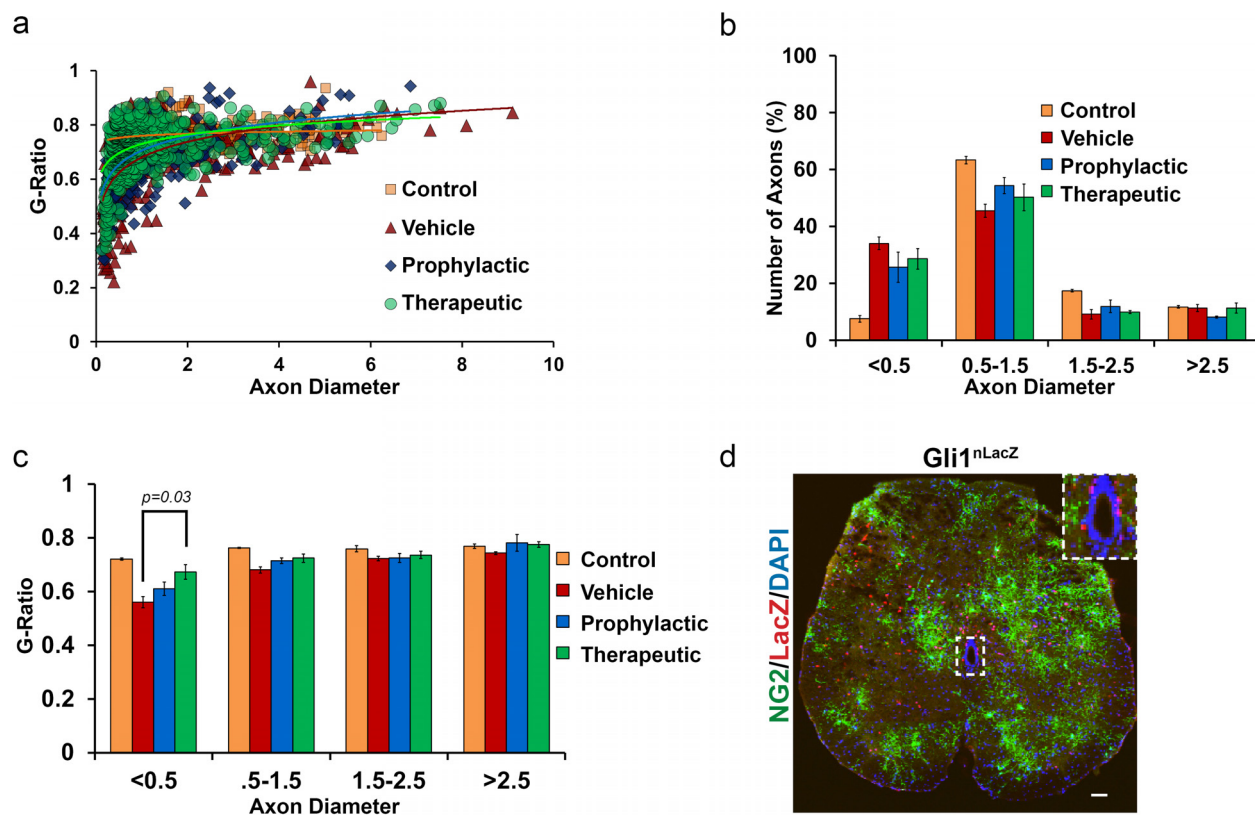
Extended Data Figure 8 | Pharmacological inhibition of Gli1 does not affect OPC recruitment or differentiation during remyelination. **a**, NG2^{CE/+} mice were treated with two doses of intraperitoneal tamoxifen to sparsely label OPCs and analysed at 2 weeks of recovery from cuprizone. Scale bar, 50 μ m. **b**, Numbers (~ 40 GFP⁺ cells per field) and **c**, proportions of GFP-labelled

OPCs (PDGFR- α) and mature oligodendrocytes (CC1) were similar in the GANT61- versus the vehicle-treated mice in the CC, indicating GANT61 does not alter OPC remyelination. Scale bar, 50 μ m. $N = 3$ mice per group. Data are mean \pm s.e.m., Student's t -test.



Extended Data Figure 9 | Gli1 is not expressed by immune cells of spleen, thymus and liver of healthy mice. a–d, Cells from the spleen, liver and thymus of tamoxifen-treated *Gli1*^{Het} (*Gli1*^{CE/+}) and *Gli1*^{Null} (*Gli1*^{CE/nLacZ}) mice were analysed by flow cytometry for GFP expression. Wild-type mice did not express GFP and were used as controls. Representative flow cytometry scatter plots

showing absence of GFP expression in CD45⁺/CD3⁺/CD4⁺/CD19⁺ T and B cells (a), CD45⁺/CD3⁺/CD19⁺/B220⁺/MPDCA⁺ plasmacytoid dendritic cells (b), CD45⁺/CD3⁺/CD19⁺/B220⁺/CD11b⁺ macrophages, monocytes, dendritic and natural killer (NK) cells (c) and CD45⁺/CD3⁺/CD19⁺/B220⁺/CD11b⁺ cells (d) in *Gli1*^{CE/+} and *Gli1*^{CE/nLacZ} mice. *N* = 3 mice per genotype.



Extended Data Figure 10 | Effects of GANT61 on spinal cord axons in the PLP-induced EAE model. **a**, Scatter plot of G ratios with respect to axonal diameters ($n = 500$ axons in three mice per group, exponential trend line). **b**, Analysis of electron microscopy images showing the relative proportion of axons binned by their diameters in the four groups. **c**, Analysis of electron microscopy images, indicating the G ratios of axons relative to their diameters

in the four groups ($N = 500$ axons in three mice per group; data are mean \pm s.e.m., Student's t -test). **d**, Immunofluorescence image of a spinal cord section from *Gli1^{nLacZ/+}* mice shows that LacZ is not expressed by NG2⁺ OPCs. The inset shows expression of LacZ in the germinal zone around the central canal. $N = 3$ mice. Scale bar, 50 μ m.

Alternative transcription initiation leads to expression of a novel *ALK* isoform in cancer

Thomas Wiesner^{1,2}, William Lee^{3,4}, Anna C. Obenaus⁵, Leili Ran¹, Rajmohan Murali^{6,7}, Qi Fan Zhang¹, Elissa W. P. Wong¹, Wenhao Hu¹, Sasinya N. Scott^{1,6}, Ronak H. Shah^{1,6}, Iñigo Landa¹, Julia Button^{1,†}, Nathalie Lailier⁷, Andrea Sboner^{8,9,10}, Dong Gao¹, Devan A. Murphy¹, Zhen Cao¹, Shipra Shukla¹, Travis J. Hollmann⁶, Lu Wang⁶, Laetitia Borsu⁶, Taha Merghoub¹¹, Gary K. Schwartz¹², Michael A. Postow^{13,14}, Charlotte E. Ariyan¹⁵, James A. Fagin^{1,13,14}, Deyou Zheng^{16,17,18}, Marc Ladanyi^{1,6}, Klaus J. Busam⁶, Michael F. Berger^{1,6,7}, Yu Chen^{1,13,14} & Ping Chi^{1,13,14}

Activation of oncogenes by mechanisms other than genetic aberrations such as mutations, translocations, or amplifications is largely undefined. Here we report a novel isoform of the anaplastic lymphoma kinase (ALK) that is expressed in ~11% of melanomas and sporadically in other human cancer types, but not in normal tissues. The novel *ALK* transcript initiates from a *de novo* alternative transcription initiation (ATI) site in *ALK* intron 19, and was termed *ALK*^{ATI}. In *ALK*^{ATI}-expressing tumours, the ATI site is enriched for H3K4me3 and RNA polymerase II, chromatin marks characteristic of active transcription initiation sites¹. *ALK*^{ATI} is expressed from both *ALK* alleles, and no recurrent genetic aberrations are found at the *ALK* locus, indicating that the transcriptional activation is independent of genetic aberrations at the *ALK* locus. The *ALK*^{ATI} transcript encodes three proteins with molecular weights of 61.1, 60.8 and 58.7 kilodaltons, consisting primarily of the intracellular tyrosine kinase domain. *ALK*^{ATI} stimulates multiple oncogenic signalling pathways, drives growth-factor-independent cell proliferation *in vitro*, and promotes tumorigenesis *in vivo* in mouse models. *ALK* inhibitors can suppress the kinase activity of *ALK*^{ATI}, suggesting that patients with *ALK*^{ATI}-expressing tumours may benefit from *ALK* inhibitors. Our findings suggest a novel mechanism of oncogene activation in cancer through *de novo* alternative transcription initiation.

To identify novel mechanisms of oncogene activation, we performed transcriptome analyses (RNA sequencing (RNA-seq)) of metastatic melanoma and thyroid carcinoma. We used an algorithm² to investigate the differential expression of exons and focused our analysis on receptor tyrosine kinases with high expression of the kinase domain. In two melanoma (MM-15, MM-74) and one anaplastic thyroid carcinoma (ATC-28) samples, we identified a novel *ALK* transcript, which contained the *ALK* exons 20–29 preceded by ~400 base pairs (bp) of intron 19, but not exons 1–19. The novel *ALK* transcript was distinct from wild-type *ALK*, which contains all exons, and from *ALK* translocations, which usually encompass exons 20–29 with little intronic expression due to preserved splice sites (Fig. 1a and Extended Data Fig. 1a–c). We confirmed the presence of the novel *ALK* transcript with a northern blot (Extended Data Fig. 2a, b).

The RNA-seq profile of the novel *ALK* transcript suggested an alternative transcription initiation site in intron 19, and we termed the novel transcript *ALK*^{ATI}. We performed 5'-rapid amplification

of cDNA ends (5'-RACE) and mapped the ATI site to a 25 bp region in intron 19 (Fig. 1b, Extended Data Fig. 2c–e and Supplementary Table 1). ChIP-seq and ChIP-qPCR showed that only *ALK*^{ATI}-expressing tumours, but not controls, had significant enrichment of histone H3K4me3 and RNA polymerase II (RNAPol II) at the ATI site, which are typical of active promoters¹ (Fig. 1c and Extended Data Fig. 3a, b). These data suggest that *ALK*^{ATI} originates from a newly established bona fide ATI site associated with characteristic chromatin alterations.

To determine the prevalence of *ALK*^{ATI} expression, we screened more than 5,000 samples from 15 different cancer types in the TCGA RNA-seq data set. *ALK*^{ATI} was expressed in ~11% of melanoma and sporadically in other cancer types (Extended Data Table 1). We found no *ALK*^{ATI} expression in more than 1,600 samples from 43 different normal tissues in the Genotype-Tissue Expression (GTEx) RNA-seq data set³, indicating that *ALK*^{ATI} is primarily expressed in a subset of cancer. To accurately distinguish and quantify the expression of *ALK*^{ATI}, wild-type *ALK*, and translocated *ALK* in clinical specimens, we developed a NanoString nCounter assay⁴ with probes in *ALK* exons 1–19, intron 19, and exons 20–29, and identified additional *ALK*^{ATI}-expressing tumours derived from formalin-fixed paraffin-embedded clinical specimens (Fig. 1d).

To determine whether somatic genomic aberrations at the *ALK* locus contribute to the establishment of the *de novo* ATI site, we performed comprehensive genetic analyses including interphase fluorescence *in situ* hybridization (FISH), array comparative genomic hybridization (aCGH), whole-genome sequencing, and ultra-deep sequencing of the *ALK* locus, but found no genomic aberrations that could account for the *de novo* expression of *ALK*^{ATI} (Extended Data Figs 4a–d and 5a–c, and Supplementary Tables 2–4). Reasoning that local genomic aberrations are usually *cis*-acting and only alter the expression of the affected allele⁵, we analysed the single nucleotide variants (SNVs) in the DNA-, RNA-, and ChIP-seq data. In all three data sets, we found similar allelic SNV frequencies, indicating that H3K4me3 decorates both *ALK* alleles and that both *ALK* alleles are actively transcribed (Fig. 1e). These data suggest that the transcriptional activation of *ALK*^{ATI} is independent of genomic aberrations at the *ALK* locus, and that alteration of *trans*-acting elements, such as transcription factors or chromatin modifiers, may contribute to *ALK*^{ATI} expression.

¹Human Oncology and Pathogenesis Program, Memorial Sloan Kettering Cancer Center, New York 10065, USA. ²Department of Dermatology, Medical University of Graz, 8010 Graz, Austria.

³Computational Biology Program, Memorial Sloan Kettering Cancer Center, New York 10065, USA. ⁴Department of Radiation Oncology, Memorial Sloan Kettering Cancer Center, New York 10065, USA.

⁵Cancer Biology and Genetics Program, Memorial Sloan Kettering Cancer Center, New York 10065, USA. ⁶Department of Pathology, Memorial Sloan Kettering Cancer Center, New York 10065, USA. ⁷Marie-Josée and Henry R. Kravis Center for Molecular Oncology, Memorial Sloan Kettering Cancer Center, New York 10065, USA. ⁸Department of Pathology and Laboratory Medicine, Weill Cornell Medical College/New York Presbyterian Hospital, New York 10065, USA. ⁹Institute for Computational Biomedicine, Weill Cornell Medical College/New York Presbyterian Hospital, New York 10065, USA. ¹⁰Institute for Precision Medicine, Weill Cornell Medical College/New York Presbyterian Hospital, New York, USA. ¹¹Immunology Program, Memorial Sloan Kettering Cancer Center 10065, New York, USA.

¹²Herbert Irving Comprehensive Cancer Center, Columbia University Cancer Center, New York 10032, USA. ¹³Department of Medicine, Memorial Sloan Kettering Cancer Center, New York 10065, USA.

¹⁴Department of Medicine, Weill Cornell Medical College, New York 10065, USA. ¹⁵Department of Surgery, Memorial Sloan Kettering Cancer Center, New York 10065, USA. ¹⁶Department of Neurology, Albert Einstein College of Medicine, New York 10461, USA. ¹⁷Department of Genetics, Albert Einstein College of Medicine, New York 10461, USA. ¹⁸Department of Neuroscience, Albert Einstein College of Medicine, New York 10461, USA. [†]Present address: Johns Hopkins School of Medicine, Baltimore, Maryland 21205-2196, USA.

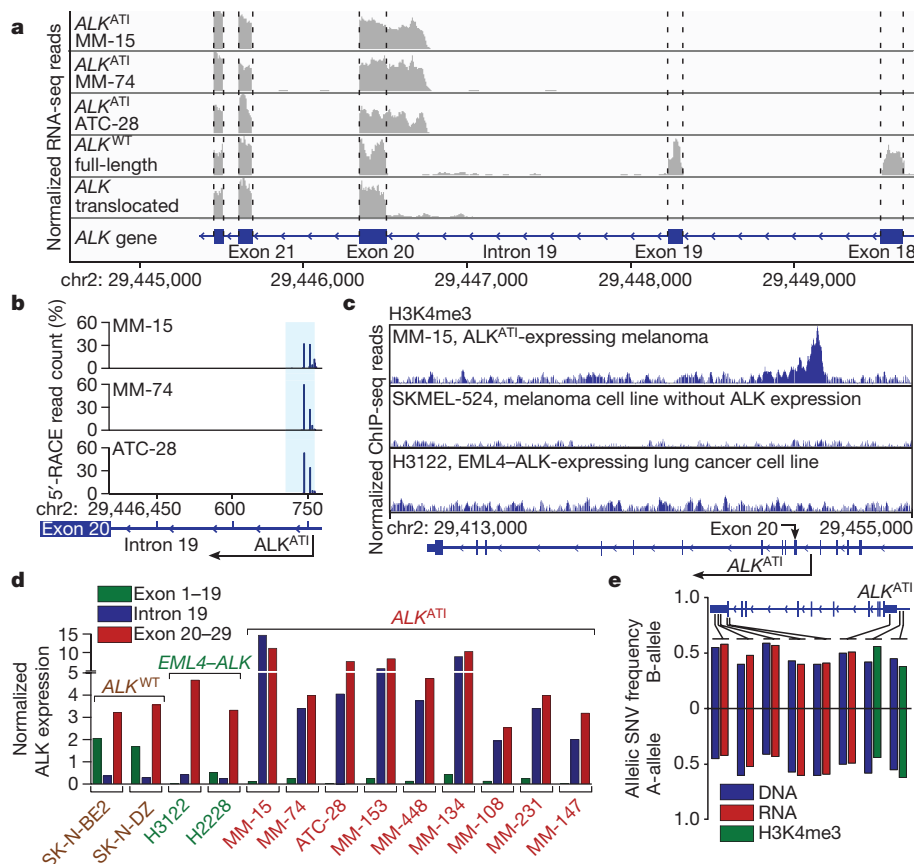


Figure 1 | Alternative transcription initiation (ATI) results in a novel *ALK* transcript. **a**, Distribution of RNA-seq reads of *ALK* variant transcripts: *ALK*^{ATI} RNA-seq reads align to both *ALK* intron 19 and exons 20–29; full-length, wild-type *ALK* (*ALK*^{WT}) RNA-seq reads align to all *ALK* exons, but not to the introns; translocated *ALK* RNA-seq reads align only to *ALK* exons 20–29. **b**, Mapping of the ATI sites of *ALK*^{ATI} to a 25 bp region in *ALK* intron 19 (hg19 chr2:29,446,768–29,446,744; blue shaded area). **c**, ChIP-seq profile of

H3K4me3 at the ATI site. **d**, Quantitative mRNA profiling of different *ALK* variants using Nanostring nCounter: 2 wild-type *ALK*-expressing neuroblastoma cell lines (SK-N-BE2 and SK-N-DZ), 2 *EML4-ALK*-expressing lung cancer cell lines (H3122 and H2228), 9 *ALK*^{ATI}-expressing tumours (8 melanomas (MM) and 1 anaplastic thyroid carcinoma, ATC-28). **e**, Similar SNV frequencies in DNA-seq, RNA-seq, and ChIP-seq (H3K4me3) data indicate that *ALK*^{ATI} is biallelically expressed.

The ATI region contains transposable elements, including a long terminal repeat (LTR) in *ALK* intron 19 and a long interspersed nuclear element (LINE) in intron 18, both of which can regulate transcription⁶ (Extended Data Fig. 6a). To evaluate whether CpG methylation of these elements might be associated with *ALK*^{ATI} expression, we performed bisulfite sequencing. Compared to the controls, the *ALK*^{ATI}-expressing samples showed lower CpG methylation in regions flanking the ATI site, including the LINE (Extended Data Fig. 6b–d). The LTR contained only few CpG sites with low methylation levels in all samples. As expected, we found H3K27ac, a histone mark characteristic of active promoters and enhancers, enriched at the ATI site of *ALK*^{ATI}-expressing tumour samples. Surprisingly, H3K27ac was also enriched in *ALK*^{ATI}-negative melanoma samples, but not in the control lung cancer cell lines or the 17 non-melanoma cell lines analysed by the ENCODE consortium⁷ (Extended Data Fig. 6e, f). By integrating ChIP, DNase I hypersensitivity, and 5'-RACE data, we defined the proximal *cis*-regulatory region on chromosome 2 as chr2:29,445,000–29,447,100 and computationally determined the potential transcription factor binding motifs⁸ (Supplementary Table 5). To test whether the LTR could function as a promoter, we used a luciferase reporter assay and found that, in contrast to lung cancer cell lines, melanoma cell lines showed low but consistent luciferase activity (Extended Data Fig. 6g). These data suggest that melanomas with H3K27ac enrichment at the ATI site might be primed to express *ALK*^{ATI}.

The *ALK*^{ATI} transcript has three predicted in-frame start codons (ATGs), resulting in proteins with molecular weights of 61.1, 60.8, and 58.7 kilodaltons (kDa). All three proteins maintain the intracellular

tyrosine kinase domain, but lack the extracellular and transmembrane domains of wild-type *ALK* (Fig. 2a). Immunoblots of two neuroblastoma cell lines (SK-N-BE2, SK-N-DZ) expressing wild-type *ALK* and two lung cancer cell lines (H3122 and H2228) expressing two distinct variants of the *EML4-ALK* gene fusion showed bands at the expected sizes. *ALK*^{ATI}-expressing tumours revealed a double band at ~60 kDa, suggesting that *ALK*^{ATI} is translated from more than one start codon (Fig. 2b). To confirm our prediction, we mutated the three start codons individually or in combination, and expressed them in 293T cells. Immunoblots revealed that each mutated form of *ALK*^{ATI} no longer produced the corresponding protein band, indicating that all three start codons in *ALK*^{ATI} are functional and give rise to three distinct proteins (Fig. 2c).

ALK^{ATI} proteins were phosphorylated in tumours with endogenous *ALK*^{ATI} expression and in cells with exogenous *ALK*^{ATI} expression (Fig. 2b, c), indicating that *ALK*^{ATI} is active. *ALK* activity was confirmed by an *in vitro* kinase assay (Extended Data Fig. 7a). A kinase-dead *ALK*^{ATI} (*ALK*^{ATI-KD}), in which a lysine in the ATP-binding site of the kinase domain was replaced by a methionine⁹, was not phosphorylated or active. Reasoning that *ALK*^{ATI} may auto-activate by forming homodimers similar to other receptor tyrosine kinases¹⁰, we tested the ability of self-interaction using co-immunoprecipitation with V5- and HA-tagged *ALK*^{ATI} proteins. The V5-*ALK*^{ATI} readily co-immunoprecipitated with the HA-*ALK*^{ATI} and vice versa, indicating that *ALK*^{ATI} can self-interact, resulting in auto-phosphorylation and kinase activity (Fig. 2d). Using immunofluorescence, we detected *ALK*^{ATI} in both the nucleus and the cytoplasm, whereas *ALK* with the

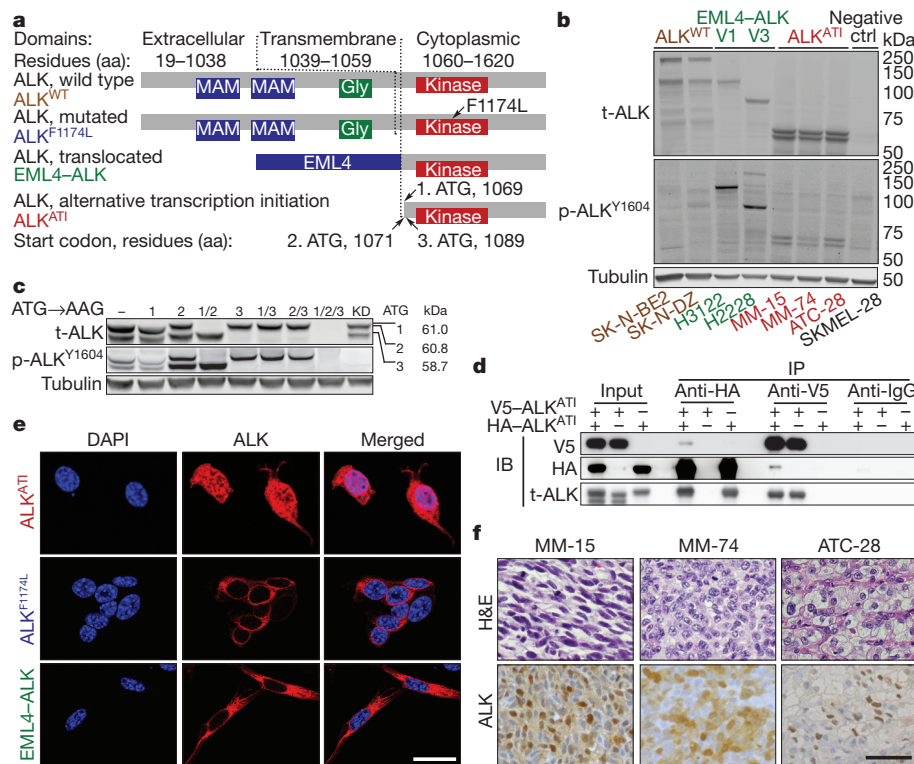


Figure 2 | The ALK^{AT1} transcript encodes three ALK proteins mainly containing the kinase domain. **a**, Illustration of ALK protein isoforms. MAM, meprin, A-5 protein, and receptor protein-tyrosine phosphatase mu; Gly, glycine-rich region. **b**, Immunoblots of total (t) and phosphorylated (p) ALK in two ALK -expressing neuroblastoma cell lines, two $EML4-ALK$ variant-expressing lung cancer cell lines, three ALK^{AT1} -expressing tumours, and a negative control melanoma cell line. **c**, Immunoblots of 293T cells with transient expression of ALK^{AT1} , in which the three predicted start codons were

mutated from ATG to AAG, individually or in combination as indicated. KD, kinase-dead. **d**, Co-immunoprecipitation (IP) and immunoblots (IB) in 293T cells expressing V5-tagged ALK^{AT1} (V5- ALK^{AT1}), HA-tagged ALK^{AT1} (HA- ALK^{AT1}), or both. **e**, ALK immunofluorescence in NIH-3T3 cells expressing the indicated ALK isoforms. Scale bar, 25 μ m. **f**, Haematoxylin and eosin (H&E) staining and ALK immunohistochemistry in ALK^{AT1} -expressing human tumour samples. Scale bar, 50 μ m. See Supplementary Fig. 1 for uncropped blots from **b–d**.

F1174L mutation (ALK^{F1174L}) and $EML4-ALK$ were found mainly in the cytoplasm and/or at the cell membrane (Fig. 2e). ALK immunohistochemistry in clinical samples confirmed the nuclear and cytoplasmic localization of ALK^{AT1} , suggesting that detection of nuclear ALK expression by immunohistochemistry could be used as a simple biomarker to identify ALK^{AT1} -expressing tumours (Fig. 2f and Extended Data Fig. 7b).

To establish the functional consequences of ALK^{AT1} expression, we stably transduced Ba/F3, NIH-3T3 and melan-a cells with ALK^{AT1} , negative controls (ALK^{AT1-KD} and empty vector), and positive controls

(oncogenic ALK variants ALK^{F1174L} , $EML4-ALK$, and wild type). In interleukin 3 (IL-3)-dependent Ba/F3 cells, expression of ALK^{AT1} and the positive controls, but not the negative controls, led to IL-3-independent proliferation (Fig. 3a). We confirmed that ALK^{AT1} was expressed at similar levels in the transformed ALK^{AT1} -Ba/F3 cells to those in human tumours, and that all ALK isoforms were phosphorylated when expressed at levels required for IL-3-independent growth (Fig. 3b). In competition assays, only Ba/F3 cells expressing green fluorescent protein (GFP) co-expressed from the ALK expression vectors were growing under IL-3-independent growth conditions,

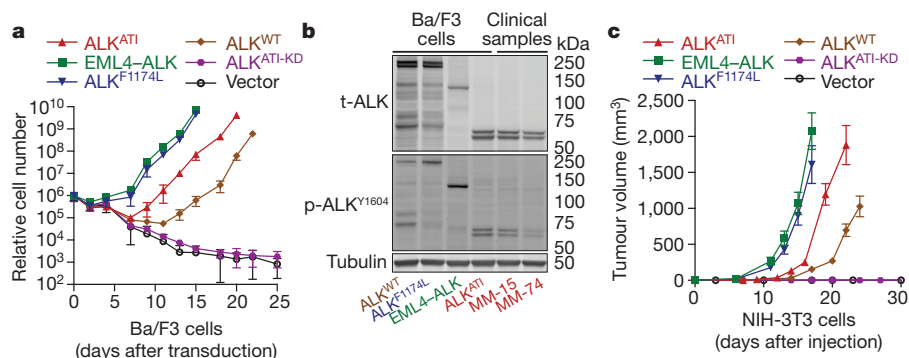


Figure 3 | ALK^{AT1} promotes growth-factor-independent proliferation *in vitro* and tumorigenesis *in vivo*. **a**, Growth curves of Ba/F3 cells stably expressing the indicated ALK isoforms in the absence of IL-3. Error bars, mean \pm s.d.; $n = 8$, pooled data from 2 experiments with 4 technical replicates each. **b**, ALK immunoblots of previously transformed (capable of IL-3-

independent growth) Ba/F3 cells with exogenous expression of the indicated ALK isoforms, and of tumours with endogenous ALK^{AT1} expression. See Supplementary Fig. 1 for uncropped blots. **c**, Tumour growth of NIH-3T3 cells stably expressing the indicated ALK isoforms. Error bars, mean \pm s.e.m.; $n = 10$ tumours, see also Source Data associated with this figure.

indicating that the Ba/F3 cell transformation was driven by expression of the *ALK* variants (Extended Data Fig. 7c). Consistently, *ALK*^{AT1}-expressing NIH-3T3 and melanoma cells efficiently induced tumour growth in mice with severe combined immunodeficiency (SCID) (Fig. 3c and Extended Data Fig. 7d–f).

All cells expressing *ALK* variants (*ALK*^{AT1}, *ALK*^{F1174L}, *EML4-ALK*, and wild-type *ALK*) were able to establish growth-factor-independent proliferation and tumorigenesis with similar growth rates once the tumours were established. The observed oncogenic capacity of wild-type *ALK* is consistent with previous reports that high endogenous expression or genomic amplification of *ALK* drives oncogenesis and confers sensitivity to *ALK* inhibitors in neuroblastomas^{11–16}. To explore the functional consequences of *ALK*^{AT1} expression further, we stably transduced NIH-3T3 cells with either a low or high titre of *ALK*^{AT1}, resulting in cells expressing *ALK*^{AT1} either at low or at high levels. We found that a further increase in *ALK*^{AT1} expression levels did not accelerate tumour graft establishment and growth, indicating that

ALK^{AT1} can drive tumorigenesis once a threshold of expression is reached (Extended Data Fig. 7g–i).

To examine the therapeutic responses to pharmacological *ALK* inhibition, we treated Ba/F3 cells stably expressing various *ALK* isoforms with three different *ALK* inhibitors (crizotinib, ceritinib, and TAE-684). All three *ALK* inhibitors effectively inhibited IL-3-independent growth of the transformed Ba/F3 cells, whereas they had no effect on growth in the presence of IL-3 (Fig. 4a and Extended Data Fig. 8a, b). Crizotinib inhibited *ALK*^{AT1} phosphorylation and downstream signalling in a concentration-dependent manner, further corroborating that *ALK*^{AT1} is activated through auto-phosphorylation (Fig. 4b and Extended Data Fig. 8c–e). Crizotinib treatment induced regression of *ALK*^{AT1}-driven NIH-3T3 tumours *in vivo*, and immunohistochemistry of explanted tumours confirmed reduced cell proliferation, increased apoptosis, and inhibition of several oncogenic signalling pathways (Fig. 4c–e and Extended Data Fig. 9a–f).

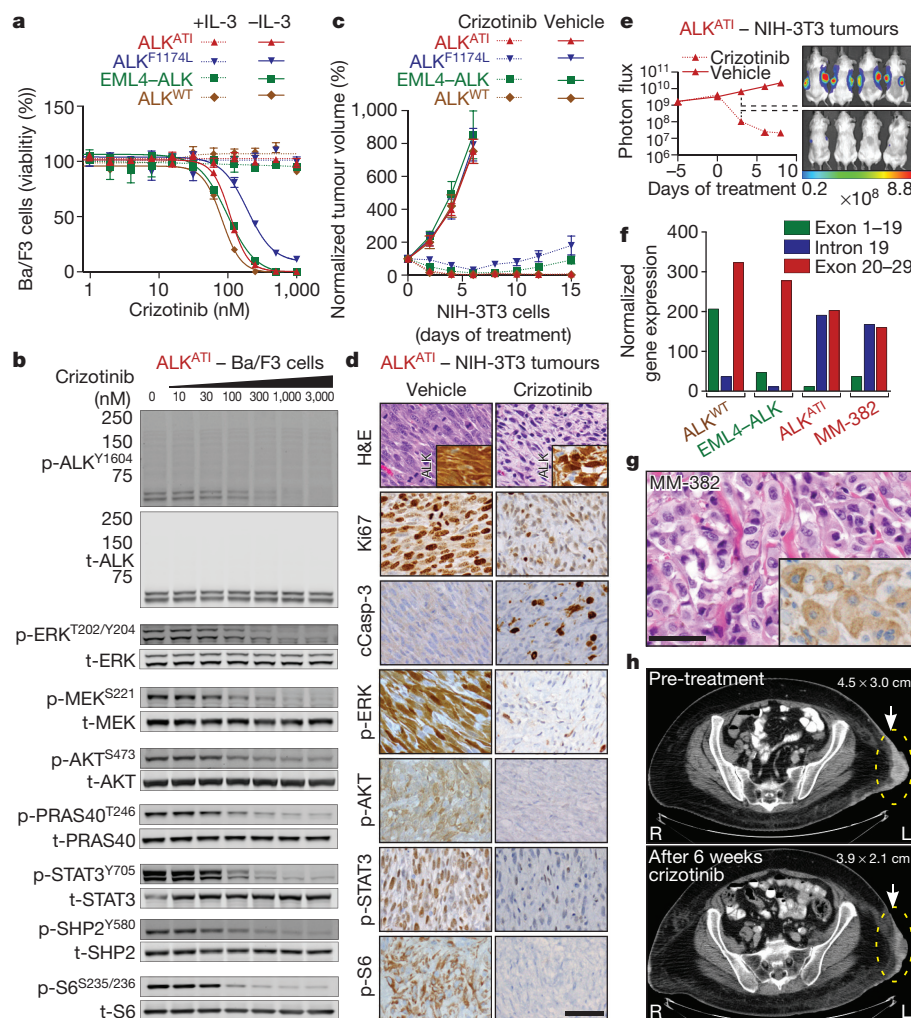


Figure 4 | *ALK*^{AT1} expression confers sensitivity to *ALK* inhibitors *in vitro* and *in vivo*. **a**, Dose–response curves to crizotinib of Ba/F3 cells expressing the indicated *ALK* isoforms in the presence or absence of IL-3. Error bars, mean \pm s.e.m.; $n = 3$ biological replicates. **b**, Representative immunoblots of *ALK*^{AT1}-expressing Ba/F3 cells treated with increasing concentrations of crizotinib for 2 h. See Supplementary Fig. 1 for uncropped blots. **c**, Normalized tumour volume in mice implanted with NIH-3T3 cells expressing the indicated *ALK* isoforms and treated with vehicle ($n = 8$ tumours) or crizotinib ($n = 10$ tumours). Error bars, mean \pm s.e.m.; see also Source Data. **d**, H&E staining and immunohistochemistry of explanted *ALK*^{AT1}-expressing tumours 48 h after

first crizotinib treatment. **e**, Normalized bioluminescence signal of *ALK*^{AT1}-expressing, luciferase-labelled NIH-3T3 tumours treated with vehicle or crizotinib. Error bars, mean \pm s.e.m.; $n = 8$ tumours; see also Source Data associated with this figure. **f**, Quantitative mRNA *ALK* profiling of a metastatic melanoma (MM-382) compared to wild-type *ALK*, *EML4-ALK*, or *ALK*^{AT1} using Nanostring nCounter. **g**, H&E staining and ALK immunohistochemistry (inset) of MM-382. Scale bars in **d** and **g**, 50 μ m. **h**, Computed tomography images of a subcutaneous melanoma metastasis from patient 1 (MM-382) before and after crizotinib treatment.

On the basis of our pre-clinical data, we identified a patient with ALK^{ATI} -expressing metastatic melanoma (Fig. 4f, g). A clinical sequencing assay¹⁷, which evaluates 341 cancer-related genes for genomic aberrations, and FISH to assess *ALK* rearrangements and *MET* amplifications, revealed deletions of *CDKN2A* and *PTEN* (Extended Data Fig. 9g–i). The patient had previously progressed on a combination of ipilimumab and nivolumab immunotherapy in a clinical trial, followed by palliative radiation and dacarbazine chemotherapy. Subsequent treatment with crizotinib resulted in marked symptomatic improvement and tumour shrinkage within 6 weeks of therapy (Fig. 4h).

Taken together, we have identified a novel *ALK* transcript, ALK^{ATI} , which arises independently of genomic aberrations at the *ALK* locus through alternative transcription initiation. ALK^{ATI} -driven tumours are sensitive to *ALK* inhibitors, suggesting that patients harbouring such tumours may benefit from *ALK* inhibitor therapy. Importantly, we have discovered alternative transcription initiation as a novel mechanism for oncogene activation. Additional oncogenes may be activated via similar mechanisms in other human malignancies, and their identification may provide new insights into oncogenesis and opportunities for therapeutic intervention.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 20 February 2014; accepted 28 July 2015.

Published online 7 October 2015.

- Kouzarides, T. Chromatin modifications and their function. *Cell* **128**, 693–705 (2007).
- Anders, S., Reyes, A. & Huber, W. Detecting differential usage of exons from RNA-seq data. *Genome Res.* **22**, 2008–2017 (2012).
- GTEX Consortium. The Genotype-Tissue Expression (GTEx) project. *Nature Genet.* **45**, 580–585 (2013).
- Reis, P. P. *et al.* mRNA transcript quantification in archival samples using multiplexed, color-coded probes. *BMC Biotechnol.* **11**, 46 (2011).
- Northcott, P. A. *et al.* Enhancer hijacking activates *GFI1* family oncogenes in medulloblastoma. *Nature* **511**, 428–434 (2014).
- Xie, M. *et al.* DNA hypomethylation within specific transposable element families associates with tissue-specific enhancer landscape. *Nature Genet.* **45**, 836–841 (2013).
- Calo, E. & Wysocka, J. Modification of enhancer chromatin: what, how, and why? *Mol. Cell* **49**, 825–837 (2013).
- Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif. *Bioinformatics* **27**, 1017–1018 (2011).
- Soda, M. *et al.* Identification of the transforming *EML4-ALK* fusion gene in non-small-cell lung cancer. *Nature* **448**, 561–566 (2007).
- Lemmon, M. A. & Schlessinger, J. Cell signaling by receptor tyrosine kinases. *Cell* **141**, 1117–1134 (2010).
- Bresler, S. C. *et al.* Differential inhibitor sensitivity of anaplastic lymphoma kinase variants found in neuroblastoma. *Sci. Transl. Med.* **3**, 108ra114 (2011).
- Mossé, Y. P. *et al.* Identification of *ALK* as a major familial neuroblastoma predisposition gene. *Nature* **455**, 930–935 (2008).
- Janoueix-Lerosey, I. *et al.* Somatic and germline activating mutations of the *ALK* kinase receptor in neuroblastoma. *Nature* **455**, 967–970 (2008).
- Passoni, L. *et al.* Mutation-independent anaplastic lymphoma kinase overexpression in poor prognosis neuroblastoma patients. *Cancer Res.* **69**, 7338–7346 (2009).
- Montavon, G. *et al.* Wild-type *ALK* and activating *ALK-R1275Q* and *ALK-F1174L* mutations upregulate *Myc* and initiate tumor formation in murine neural crest progenitor cells. *Oncotarget* **5**, 4452–4466 (2014).
- Schulte, J. H. *et al.* High *ALK* receptor tyrosine kinase expression supersedes *ALK* mutation as a determining factor of an unfavorable phenotype in primary neuroblastoma. *Clin. Cancer Res.* **17**, 5082–5092 (2011).
- Cheng, D. T. *et al.* Memorial Sloan Kettering-Integrated Mutation Profiling of Actionable Cancer Targets (MSK-IMPACT): A hybridization capture-based next-generation sequencing clinical assay for solid tumor molecular oncology. *J. Mol. Diagn.* **17**, 251–264 (2015).

Supplementary Information is available in the online version of the paper.

Acknowledgements We thank P. Romanienko for the northern blot; W. Pao for the *EML4-ALK* construct; L. Garraway for the 501Mel and WW94 cell lines; P. Koppikar for the Ba/F3 cells; and C. Sawyers, O. Abdel-Wahab, K. Griewank, and O. Guryanova for reviewing the manuscript. Next generation sequencing, array CGH and NanoString nCounter Assay were performed at the Center for Molecular Oncology, and interphase fluorescence *in situ* hybridization in the Molecular Cytogenetics Core at Memorial Sloan Kettering Cancer Center. This work was supported by grants from the Harry J. Lloyd Trust, the Jubiläumsfonds of the Österreichische Nationalbank (15461) and a Charles H. Revson Senior Fellowship to T.W.; the National Institutes of Health (NIH) grants DP2CA174499 and K08CA151660, and the Geoffrey Beene Cancer Research Fund to P.C.; the NIH grant K08CA140946, Alfred Bressler Scholars Endowment Fund and Gerstner Young Investigator Award to Y.C.; the NIH grant P50CA172012 to J. A. F.; the NIH grant P01CA12943 to M.L., M.F.B., L.W. and L.B.; and the NIH grant P30 CA008748 to Memorial Sloan Kettering Cancer Center (Core Grant).

Author Contributions Experimental design: T.W., P.C. and Y.C. Sample collection: I.L., K.J.B., M.L., T.H., J.A.F., G.K.S., L.W., T.M. and R.M. 5'-RACE, array CGH, FISH and immunohistochemistry: T.W. Preparation of DNA and cDNA libraries and bisulfite sequencing: T.W. and S.N.S. Data analyses: T.W., W.L., M.B., R.S., A.S., N.L. and D.Z. NanoString: T.W. and L.B. Immunofluorescence: T.W. Western blots, and immunoprecipitation: T.W., Q.F.Z. and J.B. ChIP and ChIP-seq: L.R. and S.S. Generation of the expression vectors: T.W., Q.F.Z. and Z.C. Luciferase reporter assay: E.W.P.W. *In vivo* assays: T.W., A.C.O. and D.A.M. FACS: T.W. and W.H. *In vitro* kinase assay: T.W. and D.G. Review of histology and immunohistochemistry: T.W., K.B. and R.M. Patient data: M.A.P. and C.E.A. Manuscript writing: T.W., Y.C. and P.C. All authors reviewed and edited the manuscript.

Author Information The sequence of ALK^{ATI} has been deposited in the European Nucleotide Archive under the accession number LN864494. RNA-seq, ChIP-seq, DNA and bisulfite sequencing data have been deposited in the NCBI Sequence Read Archive with the accession number SRP058714. Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to Y.C. (cheny1@mskcc.org) or P.C. (chip@mskcc.org).

METHODS

No statistical methods were used to predetermine sample size.

Human tumour samples. The study was approved by the Institutional Review Boards and Ethics Committees of Memorial Sloan Kettering Cancer Center, New York and informed consent was obtained from all subjects (#12-245 and #00-144). Representative portions of excised tumours were snap-frozen in liquid nitrogen, or fixed in 4% neutral buffered formalin, embedded in paraffin, and processed using routine histological methods and stained with haematoxylin and eosin.

RNA sequencing. Total RNA was extracted from fresh-frozen tissue sections (17 metastatic melanomas and 6 thyroid carcinomas) using an RNeasy Mini Kit (Qiagen). The isolated RNA was processed using the TruSeq RNA sample Prep kit (Illumina) according to the manufacturer's protocol. The libraries were sequenced on an Illumina HiSeq 2500 platform with 50, 75, or 100 bp paired-end reads to obtain on average 40,000,000–100,000,000 reads per sample. Sequencing data was mapped to hg19, and analysed using publicly available software packages: SAMtools¹⁸, Tophat¹⁹, GATK²⁰, Picard (<http://picard.sourceforge.net>), and IGV²¹.

Screening for aberrantly expressed kinases. For initial screening of RNA-seq data, candidate receptor tyrosine kinase (RTK) genes were defined by Gene Ontology annotation GO:0004714 as found in AmiGO²². DEXSeq² was used to calculate exon level counts using RTK Ensembl Gene IDs. For each gene the ratio of reads in the first half of the gene to the second half was calculated. *ALK* was identified as the top hit.

Analysis of public data sets. RNA-seq data were downloaded from the Broad Institute GTEx Genotype-Tissue Expression Portal (<http://www.broadinstitute.org/gtex/>). Level 3 TCGA data was downloaded from the Broad Institute TCGA GDAC Firehose (<http://gdac.broadinstitute.org/>) 2013_09_23 run using exon_quantification data from `illumina-hiseq_rna-seq_v2__unc_u`. *ALK*^{ATI} candidates were identified as samples with an *ALK* expression level of RSEM ≥ 100 , ≥ 500 total reads across all *ALK* exons, and ≥ 10 -fold greater average expression in exons 20–29 compared to exons 1–19. To confirm *ALK*^{ATI} expression, candidates were manually examined in IGV²¹. ENCODE ChIP-seq data for H3K27ac, mapped to hg19 and converted to bigwig track format, was downloaded from <http://genome.ucsc.edu/ENCODE/dataMatrix/encodeChIPMatrixHuman.html>.

Promoter/motif analysis. The proximal *cis*-regulatory region, chr2:29,445,000–29,447,100, was scanned for transcription factor motifs using FIMO⁸ with default parameters against the known vertebrate transcription factor motifs in the JASPAR database²³.

5'-rapid amplification of cDNA ends (5'-RACE). We used two independent 5'-RACE techniques and three *ALK*^{ATI}-expressing tumours (MM-15, MM-74, and ATC-28) to map the ATI site and the 5'-end of the *ALK*^{ATI} transcript. We applied a tobacco-acid-pyrophosphatase 5'-RACE technique according to the manufacturer's protocol (Epicentre) using the following primers: 5'-TCATACAC ATACGATTTAGGTGACACTATAGAGCGGCCGCTGCAGGAAA-3'; 5'-CAGGTCCTGATGGAGGAGGTCTTGCCAGCAAAGCA-3'. RACE products were sequenced on an Illumina MiSeq System with a 150 bp paired-end protocol according to the manufacturer's instructions. The sequencing reads were mapped to hg19 using BWA and visualized using IGV²¹. We confirmed the continuous transcription starting in *ALK* intron 19 with an independent oligonucleotide-based 5'-RACE kit (Clontech) according to the manufacturer's protocol using the primers 5'-CTAATACGACTCACTATAGGGC-3', 5'-ACACCTGGCC TTCATACACCTCC-3'. We cloned the RACE cDNA products into plasmids (Invitrogen) and analysed them with Sanger sequencing. Two lung cancer cell lines (H3122 and H2228) with *EML4-ALK* translocations were used as controls.

Chromatin immunoprecipitation sequencing (ChIP-seq) and ChIP-qPCR. Chromatin was isolated from human tumour tissue and cell lines. Fresh-frozen human tumour tissue (MM-15, MM-74, and ATC-28) was sectioned with a microtome and cross-linked in 1% paraformaldehyde for 15 min. The cross-linked tissue samples were quenched in 125 mM glycine, washed in PBS, re-suspended in lysis buffer, ground in a Tenbroeck-style tissue grinder, and sonicated. Chromatin isolation from cell lines and immunoprecipitation was performed as previously described²⁴. Solubilized chromatin was immunoprecipitated with antibodies against H3K4me3, H3K27ac, and RNAPII (all Active Motif). Sequencing was performed on an Illumina HiSeq 2500 with 51 bp single reads. Reads were aligned to hg19 using Bowtie within the Illumina Analysis Pipeline. Peak calling was performed using MACS 1.4 comparing immunoprecipitated chromatin with input chromatin²⁵. ChIP-qPCR was performed on a ViiA 7 Real Time PCR System (Life Technologies) using Power SYBR Master Mix (Life Technology) with 3 technical replicates and 5 independent primers pairs, which are specified in Supplementary Table 6.

Ultra-deep targeted sequencing of the entire *ALK* locus and MSK-IMPACT. Targeted sequencing of the entire *ALK* locus was performed using custom hybridization capture probes tiling hg19 chr2:29,400,000–30,300,000 (Roche/

NimbleGen's SeqCap EZ). This region encompassed the entire genomic footprint of *ALK* and ~150 kb of upstream sequence. After the genomic DNA was fragmented (E220, Covaris), we prepared barcoded sequence libraries (New England Biolabs, Kapa Biosystems) and performed hybridization capture on barcoded pools. Using 250 ng of genomic DNA, we constructed libraries from 7 separate samples: 2 melanomas (MM-15 and MM-74), 1 anaplastic thyroid carcinoma (ATC-28), 2 lung cancer cell lines (H3122 and H2228) with *EML4-ALK* translocations, 1 melanoma cell line (SKMEL-28), and 1 control pool of 10 'normal' blood samples. Libraries were pooled at equimolar concentrations (100 ng per library) and used in the capture reaction as previously described²⁶. To prevent off-target hybridization, we added a pool of spike-in blocker oligonucleotides complementary to the full sequences of all barcoded adaptors. The captured libraries were sequenced on an Illumina HiSeq 2500 to generate 75 bp paired-end reads. Sequence data were de-multiplexed using CASAVA, and aligned to hg19 using BWA²⁷. Local realignment and quality score recalibration were performed using the Genome Analysis Toolkit (GATK) according to GATK best practices²⁸. We achieved mean unique target sequence coverage of $1,778 \times$ per sample (range: $1,293$ – $2,188 \times$). Sequence data were analysed to identify single nucleotide variants, small insertions/deletions (indels), and structural rearrangements. Single nucleotide variants were called using muTect²⁹ and were compared to the negative control pool (pooled 'normal' blood samples). Variants were retained if the variant allele frequency in the tumour was > 5 times than in the negative control and the frequency in the negative control was < 0.02 . Validated SNPs in the dbSNP database were filtered out. Indels were called using the SomaticIndelDetector tool in GATK²⁸ and were retained if the tumour harboured > 3 supporting reads and the frequency in the negative control was < 0.02 . DELLY was used to search for structural rearrangements³⁰. Memorial Sloan Kettering-Integrated Mutation Profiling of Actionable Cancer Targets (MSK-IMPACT) was performed as described previously¹⁷.

Bisulfite sequencing of the entire *ALK* locus. We performed custom capture of the entire *ALK* locus using custom hybridization capture probes tiling the entire genomic footprint of *ALK* (900 kb, chr2:29,400,000–30,300,000) followed by bisulfite sequencing. After fragmentation (E220, Covaris) of 3 μ g genomic DNA of each sample (MM-15, ATC-28, H3122, and SKMEL-28), libraries were prepared with the KAPA Hyper Prep Kit (Kapa Biosystems) without PCR amplification to preserve the methylation status. Of each barcoded library, 1 μ g was pooled at equimolar concentrations and captured according to the manufacturer's protocol (Roche/NimbleGen's SeqCap EZ). After washing the Dynabeads M-270 (Life Technologies), the non-biotinylated tumour/cell line DNA was dissociated from the biotinylated capture beads with 0.5 M NaOH. The single-stranded eluted DNA was used for bisulfite conversion using the EZ DNA Methylation-Gold Kit (Zymo Research) according to the manufacturer's protocol, except for the 98 °C denaturation step. After bisulfite conversion, we used the KAPA HiFi Uracil PCR polymerase (Kapa Biosystems) to amplify the library, purified the reaction with Agencourt AMPure XP beads (A63881, Beckman Coulter), and sequenced the library on an Illumina MiSeq with a 150 bp paired-end protocol according to the manufacturer's instructions. Sequence data were aligned to hg19 and analysed using Bismark³¹. We compared the methylation level at CpG sites across all samples; no methylation was detected in the CHG and CHH contexts. Methylation was first computed as the number of methylated CpG reads vs the number of total reads covering each CpG site (sites with < 10 reads were excluded). A sliding window was used to determine the mean methylation level for every 250 bp region (with at least three CpGs) near the *ALK* promoter region (chr2:29,444,000–29,452,000). Differential methylation was evaluated using a Mann–Whitney test.

Whole-genome sequencing. Whole-genome sequencing was performed at the New York Genome Center. Briefly, genomic DNA libraries were prepared from MM-15 and ATC-28 (no matched normal DNA was available) using the Illumina PCR-free kit. Libraries were sequenced on an Illumina HiSeq 2500 using the 100 bp paired-end whole-genome sequencing protocol. Sequence reads were mapped using BWA²⁷ and processed using GATK²⁸. Genome-wide analyses of mutations (HaplotypeCaller²⁸), copy number alterations (FREENC)³², and structural variations (CREST)³³ were performed. Mutations were annotated with the Ensembl Variant Effect Predictor³⁴ and filtered to remove common polymorphisms. Non-synonymous mutations along with copy number alterations and structural variations were visualized using Circos³⁵.

Allelic frequencies SNVs at the *ALK* locus. DNA-seq, RNA-seq, and ChIP-seq (H3K4me3) data were displayed in the Integrative Genomics Viewer and the allelic frequencies were compared for each single nucleotide variant (SNV).

Array CGH. Genome-wide analysis of DNA copy number changes was conducted using an oligonucleotide SurePrint G3 Human CGH Microarray (Agilent) containing 1,000,000 probes. Slides were scanned using a microarray scanner G2505B (Agilent) and analysed using Genomic Workbench (Agilent).

Interphase fluorescence *in situ* hybridization (FISH). *ALK* break-apart probes and locus-specific probes for *MET* and for centromere 7 were purchased from Abbott. The probes were hybridized on 5- μ m-thick tissue sections and the number and localization of the hybridization signals was assessed in a minimum of 100 interphase nuclei with well-delineated contours.

Northern blots. Total RNA was extracted from fresh-frozen tissue or cell lines using Qiagen's RNeasy Mini Kit (Qiagen). Up to 10 μ g RNA was used for running formaldehyde-based northern blot analysis according to the manufacturer's protocol using the RNA Ambion NorthernMax Kit (Ambion). After hybridization with a 32 P-labelled probe, consisting of *ALK* exon 20–29, the membrane was washed and visualized.

NanoString. Details of the nCounter Analysis System (NanoString Technologies) were reported previously³⁶. In brief, two sequence-specific probes were constructed for *ALK* exons 1–19, intron 19, and exons 20–29, respectively. Four control genes (*RPS13*, *RPL27*, *RPS20*, and *ACTB*) were used for normalization. The probes were complementary to a 100 bp region of the target mRNA and are listed in Supplementary Table 7. 100 ng of total RNA from each sample was hybridized, the raw data were normalized to the standard curve generated via the nCounter system, and the average value of the two probes in each target region (exons 1–19, intron 19, exons 20–29) was printed in bar charts using GraphPad Prism software 6.0. NanoString experiments were independently performed at least twice with appropriate positive and negative controls, and a representative experiment is shown.

Cell lines. The NIH-3T3 mouse embryonic fibroblast cells (#CRL-1658) and the two neuroblastoma cell lines, SK-N-DZ (#CRL-2149) and SK-N-BE2 (#CRL-2271), were obtained from the American Type Culture Collection (ATCC), and maintained in Dulbecco's modified eagle medium (DMEM). The lung carcinoma cell lines H2228 (#CRL-5935, ATCC) and H3122 (NCI, Bethesda, MD) were cultured in Roswell Park Memorial Institute medium (RPMI). The IL-3-dependent murine pro B-cell line, Ba/F3, was obtained from the 'Deutsche Sammlung von Mikroorganismen und Zellkulturen' (DSMZ) and was cultured in RPMI supplemented with 1 ng ml⁻¹ IL-3 (R&D). The melanoma cell lines A375, A2058, and Colo800 were provided by the laboratory of J. Massagué and were cultured in DMEM. The melanoma cell lines 501mel and WW94 were provided by the laboratory of L. Garraway and were cultured in DMEM. SKMEL-23, SKMEL-28, SKMEL-31, and SKMEL-524 are patient-derived melanoma cell lines established at Memorial Sloan Kettering Cancer Center and were cultured in RPMI. Melanoma cells were provided by D. Bennett (St. George's Hospital, University of London, London, UK)³⁷ and were maintained in RPMI supplemented with 200 nM 12-O-tetradecanoylphorbol-13-acetate (TPA; Cell Signaling). For retrovirus production, 293T cells were purchased (Clontech) and cultured in DMEM. All cell culture media contained 10% FBS, L-glutamine (2 mM), penicillin (100 U ml⁻¹), and streptomycin (100 μ g ml⁻¹). All cells were cultured at 37 °C in 5% CO₂ and were biochemically tested negative for mycoplasma contamination every 4–8 weeks (Lonza). After receiving the cell lines from the indicated sources, the cell lines were not further authenticated, but showed the expected genomic aberrations, such as *MLL4-ALK* translocations. The genomic aberrations have been validated by various methods, including RNA-seq or northern blot.

Plasmids. For the *ALK*^{ATT} vector, RNA from MM-15 was reverse-transcribed with anchored oligo(dT) primers into cDNA (Roche), PCR amplified with *ALK* the primers 5'-CACCATCCCATCTCCAGTCTGCTTC-3', 5'-AGAGAAGTGA GTGTGCGACC-3', and cloned into a pENTR vector (Life Technologies). The full-length *ALK* plasmid (HsCD00079531) was purchased from the DF/HCC DNA Resource Core (<http://plasmid.med.harvard.edu>) and *MLL4-ALKv1* was synthesized at GeneArt (Life Technologies). Site-directed mutagenesis was performed using QuikChange (Agilent): for kinase-dead *ALK*^{ATT} (*ALK*^{ATT-KD}), we mutated the lysine in the ATP-binding site of the kinase domain to methionine (p.K1150M referring to wild-type *ALK*) in *ALK*^{ATT}, and for *ALK*^{F1174L} a p.F1174L mutation was introduced into wild-type *ALK*. Plasmids were sub-cloned into pMIG-w backbones³⁸ (<http://www.addgene.org>), resulting in MSCV-*ALK*^{variant}-IRES-GFP constructs, which were confirmed by digestion and sequencing. To confirm the start codons, we mutated the three start codons from ATG to AAG individually or in combination as indicated. For co-immunoprecipitation, *ALK*^{ATT} was subcloned into pcDNA3.1/nV5-Dest (Life Technologies) and MSCV-N-HA-FLAG-Dest (Addgene). For bioluminescence imaging, we used a triple modality retroviral reporter plasmid (red fluorescent protein (RFP)-thymidine-kinase-luciferase)³⁹.

Stable gene expression. Retrovirus was produced in 293T cells by standard methods using ecotropic or amphotropic packaging vectors and XtremeGene 9 (Roche). We harvested the virus-containing supernatant 48, 64 and 72 h after transfection. The supernatant was pooled, filtered through a 0.45 μ m PVDF membrane, and used for transduction in the presence of polybrene (8 μ g ml⁻¹). Stable eGFP- or RFP-expressing cells were sorted with a FACSAria II (BD Biosciences).

Co-immunoprecipitation. V5-*ALK*^{ATT} and HA-*ALK*^{ATT} was transiently transfected into 293T cells using XtremeGene 9 (Roche), and after 24 h, cells were lysed in 10 mM Tris-HCl (pH 7.5), 1% Triton X-100, 150 mM NaCl, 1 mM EDTA, 1 mM DTT, 1 mM PMSF, and proteinase/phosphatase inhibitors. After incubation and centrifugation, 100 μ l supernatant was used as input, and 300 μ l for immunoprecipitation using the following antibodies: 2 μ g of anti-V5 antibody (Thermo Scientific), 10 μ l of EZview Red anti-HA Affinity Gel (Sigma), 2 μ g of anti-mouse IgG (Santa Cruz). We used 20 μ l of Protein A/G UltraLink Resin (Thermo Scientific) for immunoprecipitation. The immunoprecipitated material was eluted in 4 \times SDS loading buffer for immunoblotting. Co-immunoprecipitation was independently performed twice and a representative immunoblot is shown in Fig. 2d.

In vitro kinase assay. Stably transduced NIH-3T3 cells were grown in a 15 cm dish, washed with PBS, and lysed in 20 mM Tris (pH 8.0), 1% NP-40, 125 mM NaCl, 2.5 mM MgCl₂, and 1 mM EDTA with proteinase/phosphatase inhibitor. Lysates were incubated on ice, centrifuged, pre-cleared with 25 μ l Protein A/G UltraLink Resin (Thermo Scientific) for 30 min at 4 °C under rotation, and immunoprecipitated with 10 μ l ALK (D5F3) XP Rabbit monoclonal antibody and 25 μ l Protein A/G UltraLink Resin. After rotation for 120 min at 4 °C, the immunoprecipitated material was washed and used according to the instructions of Universal Tyrosine Kinase Assay Kit (Clontech). After the enzymatic reaction, the immunoprecipitated material was mixed with 4 \times SDS loading buffer for immunoblotting. *In vitro* kinase assays were performed in quadruplicates, independently repeated three times, and a representative experiment is shown.

Immunohistochemistry. Immunohistochemistry was performed on archival formalin-fixed paraffin-embedded tumour specimens using a standard multimer/diaminobenzidine (DAB) detection protocol on a Discovery Ultra system (Roche/Ventana) with appropriate negative and positive controls. The following antibodies (Cell Signaling Technology) were diluted in SignalStain antibody diluent (Cell Signaling Technology) as indicated: *ALK* 1:250; phospho-Akt (Ser473) 1:50; phospho-STAT3 (Tyr705) 1:400; phospho-S6 (Ser235/236) 1:400; phospho-p44/42 MAPK (Erk1/2) (Thr202/Tyr204) 1:400; cleaved caspase-3 (Asp175) 1:400. The anti-Ki67 antibody was diluted 1:600 (Abcam).

Immunofluorescence. Stably transduced NIH-3T3 cells were grown on coverslips, fixed in 4% formaldehyde, washed in PBS, and incubated in blocking solution (5% goat serum and 0.1% Triton X-100 in PBS). After blocking, cells were incubated with an *ALK* monoclonal antibody (Cell Signaling Technology) diluted 1:1,000 in blocking buffer overnight at 4 °C. After washing the cells with 0.05% Tween 20 and PBS, cells were incubated with a secondary antibody (Life Technologies) diluted 1:500 in blocking buffer for 2 h at room temperature. After washing in PBS, slides were mounted with Prolong Gold Antifade Reagent with DAPI (Cell Signaling Technology) and imaged with a Leica TCS SP5 II confocal microscope. Immunofluorescence was independently performed twice and a representative experiment is shown in Fig. 2e.

Immunoblot. Cell lysates were prepared in RIPA buffer supplemented with proteinase/phosphatase inhibitor. Proteins were resolved in NuPAGE Novex 4–12% Bis-Tris Protein Gels (Life Technologies) and transferred electrophoretically onto a nitrocellulose 0.45 μ m membrane (BioRad). Membranes were blocked for 1 h at room temperature in Odyssey Blocking Buffer (LI-COR) and were incubated overnight at 4 °C with the primary antibodies diluted at 1:1,000 in 50% Odyssey Blocking Buffer in PBS plus 0.1% Tween 20. The following primary antibodies were used (all from Cell Signaling Technologies unless stated otherwise): anti- α -tubulin (Sigma-Aldrich), anti-V5 (Thermo Scientific), anti-HA3F10 (Roche), phospho-*ALK* (Tyr1604), *ALK*, phospho-Akt (Ser473), Akt, phospho-STAT3 (Tyr705), STAT3, phospho-S6 (Ser235/236), S6, phospho-p44/42 MAPK (Erk1/2) (Thr202/Tyr204), p44/42 MAPK (Erk1/2), phospho-MEK1/2 (Ser221), MEK1/2, phospho-PRAS40 (Thr246), PRAS40 (D23C7), phospho-SHP-2 (Tyr580), SHP-2. After 4 washes of 5 min in PBS-T, membranes were incubated with secondary antibodies (IRDye 800CW goat anti-Rabbit, 1:20,000, LI-COR; IRDye 680RD goat anti-mouse, 1:20,000, LI-COR) in 50% Odyssey Blocking Buffer in PBS plus 0.1% Tween 20 for 45 min at room temperature. After 4 washes in PBS-T and a final wash with PBS, membranes were scanned with a LI-COR Odyssey CLx scanner and adjusted using LI-COR Image Studio. Immunoblots were independently performed at least twice and a representative experiment is shown in Figs 2b–d 3b, 4b and Extended Data Figs 7a, 7d, 7f, 7h, 8c–e.

Luciferase reporter assay. The long terminal repeat in *ALK* intron 19 at the AT1 site (*LTR16B2*, chr2:29,446,649–29,447,062; 414 bp) was amplified using genomic DNA from patient MM-15 and 5'-GTCCCTCATGGCTCAGCTTGT-3' and 5'-AGCACTACACAGGCCACTTC-3' primers. The PCR product (chr2:29,446,444–29,447,174; 731 bp) was cloned into pGL4.14-firefly luciferase vector (Promega). To determine the promoter activity of *LTR16B2*, 10⁵ cells were transfected in triplicates with 500 ng pGL4.14-LTR16B2 or vector alone; as internal control, 200 ng pRL-TK-Renilla luciferase reporter vector (Promega) was

co-transfected. Luciferase activity was measured using Dual-Glo Luciferase Assay System (Promega) 48 h after transfection. Promoter activity was calculated by normalizing firefly luciferase activity to the control Renilla luciferase activity and compared between pGL4.14-LTR16B2 and vector alone. The luciferase reporter assays were independently performed three times, the results were combined, and the mean \pm s.d. is shown in Extended Data Fig. 6g.

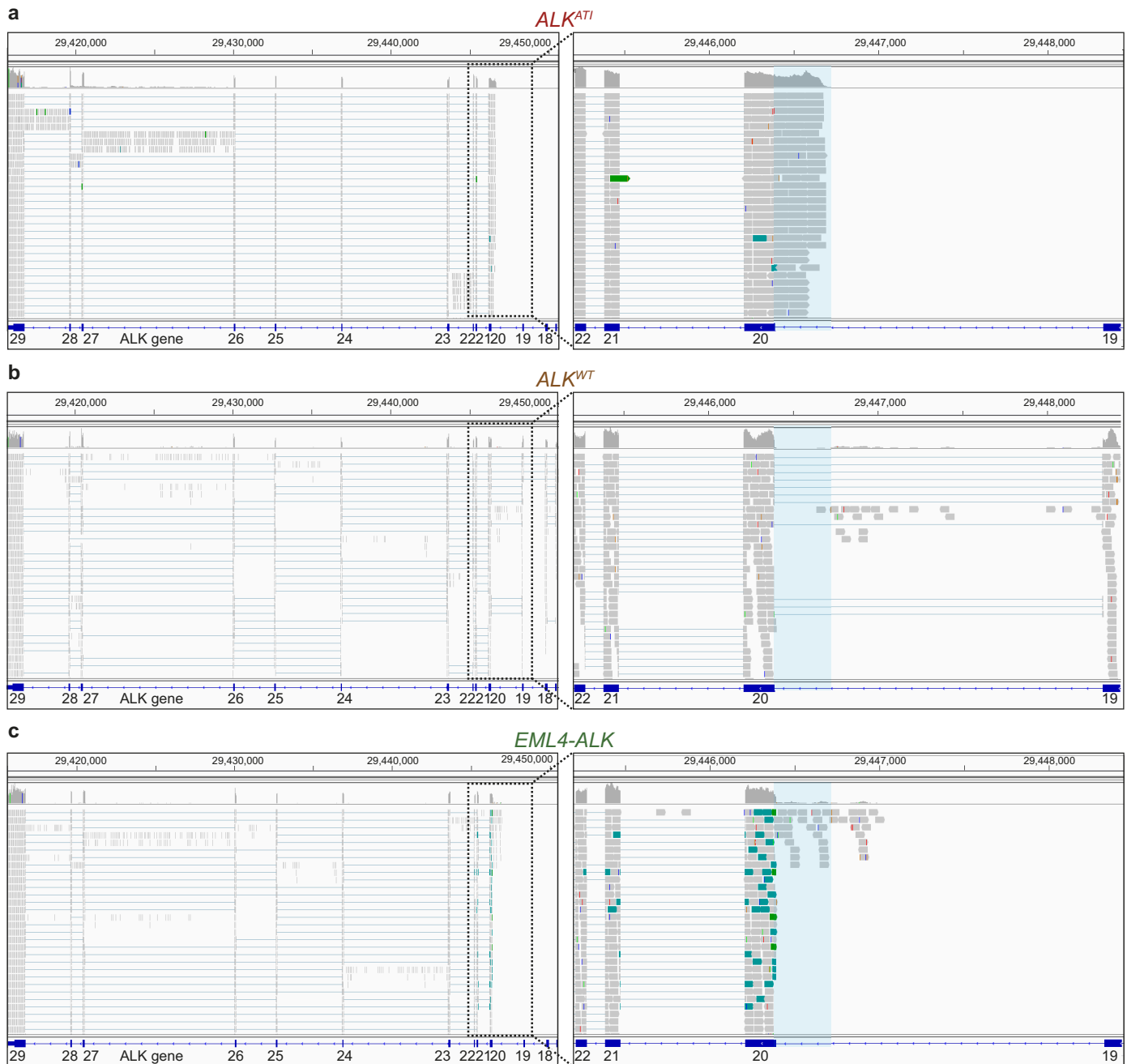
Flow cytometry and fluorescence-activated cell sorting (FACS). Flow cytometry analysis for *in vitro* transformation assays with Ba/F3 cells was performed on an LSRFortessa (BD Biosciences) at day 0 and day 14 after IL-3 withdrawal. GFP- or RFP-positive cells were sorted using the FITC (blue laser) or PE (yellow laser) channel, respectively, on a FACS Aria II configured with 5 lasers (BD Biosciences).

***In vitro* transformation and drug treatment assays.** Stably transduced Ba/F3 cells were cultured in RPMI medium supplemented with IL-3 (1 ng ml^{-1}). For the cell proliferation assay, Ba/F3 cells were transferred into IL-3 depleted RPMI medium, and cell growth was quantified in quadruplicates every 2–4 days by a luminescence assay (Promega). For cell viability assays and ALK inhibitor–dose-response curves, 2,000 Ba/F3 cells were plated in triplicates in 96-well plates with increasing concentrations of the ALK inhibitors crizotinib (LC laboratories), TAE-684 (ChemieTek), or ceritinib (ChemieTek) as indicated. All drugs were re-suspended in DMSO. The cell viability was assessed after 72 h by a luminescence assay (Promega). Results were normalized to cell growth in medium containing an equivalent concentration of DMSO. The inhibition curve was determined with GraphPad Prism 6.0 software using the ‘log(inhibitor) vs response – variable slope’ nonlinear regression model. For immunoblots, 10 million Ba/F3 cells were harvested after 2 h treatment with crizotinib, washed in ice-cold PBS, and lysed in RIPA buffer. All assays were independently performed at least twice and a representative experiment is shown.

***In vivo* tumorigenicity and drug treatment assays.** All animal experiments were performed in accordance with a protocol approved by MSKCC Institutional Animal Care and Use Committee (#11-12-029). The size for each cohort was determined based on previous experience without specific statistical methods. We re-suspended 10^6 cells of stably transfected NIH-3T3 or melan-a cells in 50 μl of 1:1 mix of PBS and Matrigel (BD Biosciences), and subcutaneously and bilaterally injected the cells into the flanks of 6–8 weeks old female CB17-SCID mice (Taconic). Mice were chosen randomly and no animals were excluded. For tumour growth assays, 5 mice were injected with parental or stably transduced cell lines and 10 tumours were assessed (expect of melan-a cells stably transfected with ALK^{F1174L}, in which 4 mice were injected and 8 tumours were assessed). Tumour sizes were measured with callipers, without blinding, every 2 to 7 days for a period of up to 100 days, and were calculated using the following formula: tumour volume = $(D \times d^2)/2$, whereby D and d refer to the long and short tumour diameter, respectively. For *in vivo* drug sensitivity studies, 9 mice were injected with the stably transduced NIH-3T3 cells expressing a luciferase reporter construct and the indicated plasmids. When the tumours reached an average size of 200–250 mm^3 , mice were randomized into a vehicle (4 mice) or treatment (5 mice) group. Mice were orally gavaged once a day with crizotinib (100 mg kg^{-1} per day) or vehicle. We performed non-blinded measurement of 8 tumours in the vehicle group and 10 tumours in the crizotinib group with callipers every 2 to 3 days, and growth curves were visualized with Prism GraphPad 6.0. In parallel, we monitored tumour growth by bioluminescence imaging of anaesthetized mice by retro-orbital

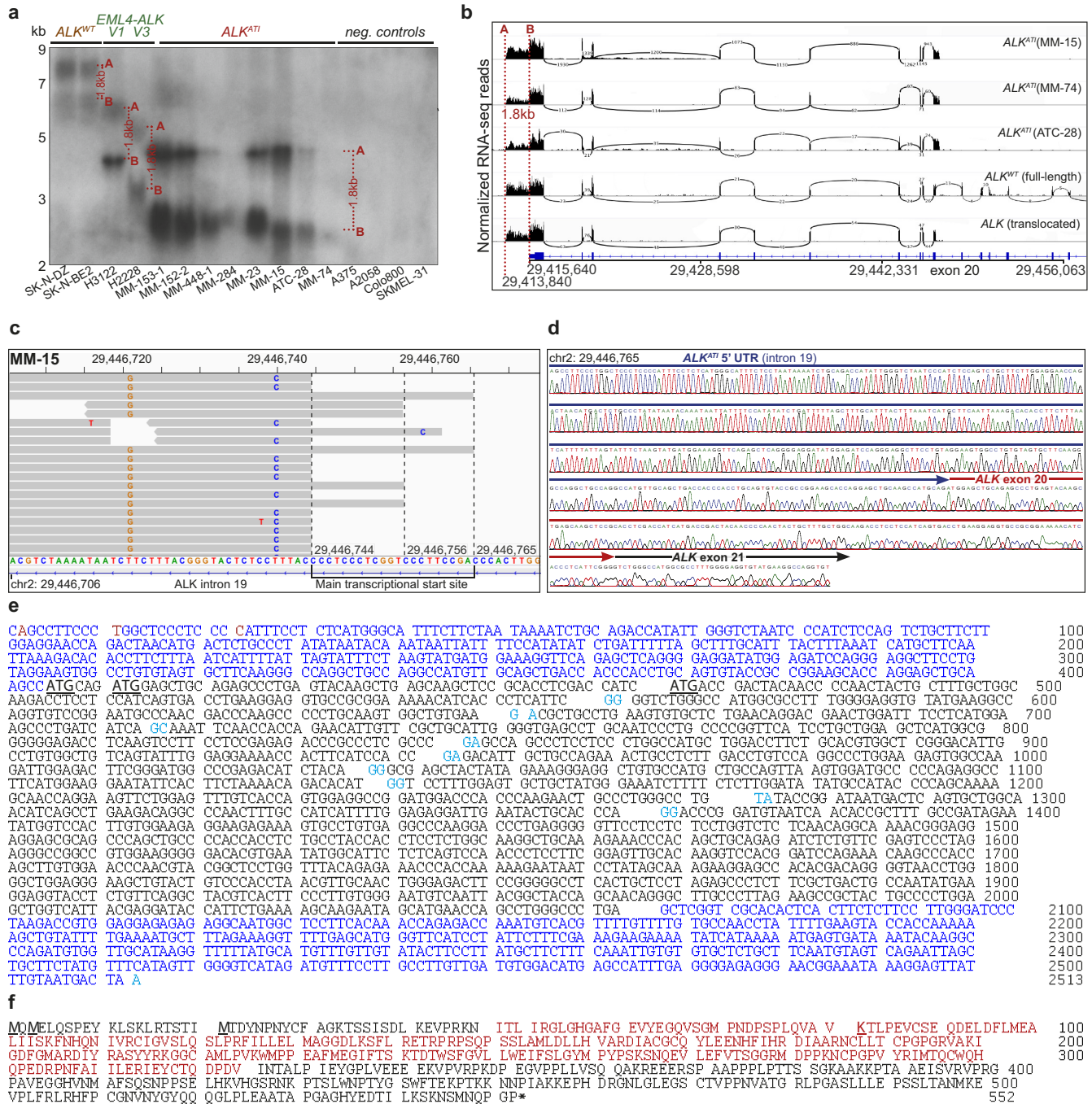
injection of D-luciferin (150 mg per kg body weight) and imaging with the IVIS Spectrum Xenogen machine (Caliper Life Sciences). Bioluminescence analysis was performed using Living Image software, version 2.50. After euthanizing the mice, tumours were explanted and either lysed in RIPA buffer (Cell Signaling Technology) or fixed overnight in 4% paraformaldehyde, washed, embedded in paraffin, and sectioned for haematoxylin and eosin (H&E) staining or immunohistochemistry. Mice experiments were performed once.

18. Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
19. Trapnell, C., Pachter, L. & Salzberg, S. L. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* **25**, 1105–1111 (2009).
20. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
21. Robinson, J. T. *et al.* Integrative genomics viewer. *Nature Biotechnol.* **29**, 24–26 (2011).
22. Carbon, S. *et al.* AmiGO: online access to ontology and annotation data. *Bioinformatics* **25**, 288–289 (2009).
23. Mathelier, A. *et al.* JASPAR 2014: an extensively expanded and updated open-access database of transcription factor binding profiles. *Nucleic Acids Res.* **42**, D142–D147 (2014).
24. Chi, P. *et al.* ETV1 is a lineage survival factor that cooperates with KIT in gastrointestinal stromal tumours. *Nature* **467**, 849–853 (2010).
25. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* **9**, R137 (2008).
26. Won, H. H., Scott, S. N., Brannon, A. R., Shah, R. H. & Berger, M. F. Detecting somatic genetic alterations in tumor specimens by exon capture and massively parallel sequencing. *J. Vis. Exp.* e50710 (2013).
27. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
28. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature Genet.* **43**, 491–498 (2011).
29. Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nature Biotechnol.* **31**, 213–219 (2013).
30. Rausch, T. *et al.* DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics* **28**, i333–i339 (2012).
31. Krueger, F. & Andrews, S. R. Bismark: a flexible aligner and methylation caller for Bisulfite-Seq applications. *Bioinformatics* **27**, 1571–1572 (2011).
32. Boeva, V. *et al.* Control-free calling of copy number alterations in deep-sequencing data using GC-content normalization. *Bioinformatics* **27**, 268–269 (2011).
33. Wang, J. *et al.* CREST maps somatic structural variation in cancer genomes with base-pair resolution. *Nature Methods* **8**, 652–654 (2011).
34. McLaren, W. *et al.* Deriving the consequences of genomic variants with the Ensembl API and SNP Effect Predictor. *Bioinformatics* **26**, 2069–2070 (2010).
35. Krzywinski, M. *et al.* Circos: an information aesthetic for comparative genomics. *Genome Res.* **19**, 1639–1645 (2009).
36. Geiss, G. K. *et al.* Direct multiplexed measurement of gene expression with color-coded probe pairs. *Nature Biotechnol.* **26**, 317–325 (2008).
37. Bennett, D. C., Cooper, P. J. & Hart, I. R. A line of non-tumorigenic mouse melanocytes, syngeneic with the B16 melanoma and requiring a tumour promoter for growth. *Int. J. Cancer* **39**, 414–418 (1987).
38. Refaelli, Y., Van Parijs, L., Alexander, S. I. & Abbas, A. K. Interferon gamma is required for activation-induced death of T lymphocytes. *J. Exp. Med.* **196**, 999–1005 (2002).
39. Ponomarev, V. *et al.* A novel triple-modality reporter gene for whole-body fluorescent, bioluminescent, and nuclear noninvasive imaging. *Eur. J. Nucl. Med. Mol. Imaging* **31**, 740–751 (2004).



Extended Data Figure 1 | Comparison of the RNA-seq profiles of various *ALK* transcripts. RNA-seq data are displayed in the Integrative Genomics Viewer (IGV). The grey bars/arrows indicate the sequencing reads. The blue lines connect sequencing reads that are aligned over the splice site of joining exons. **a.** The *ALK*^{ATI} transcript shows expression of *ALK* exons 20–29 and of ~400 bp in intron 19 (blue shaded area). No expression of exon 1–19 or intronic areas, other than in intron 19, are observed. The detailed view illustrates that the sequencing reads align continuously between exon 20 and intron 19 indicating uninterrupted transcription. The 5'-UTR of *ALK*^{ATI} (intron 19) and exon 20–29 are expressed at comparable levels. **b.** The full-

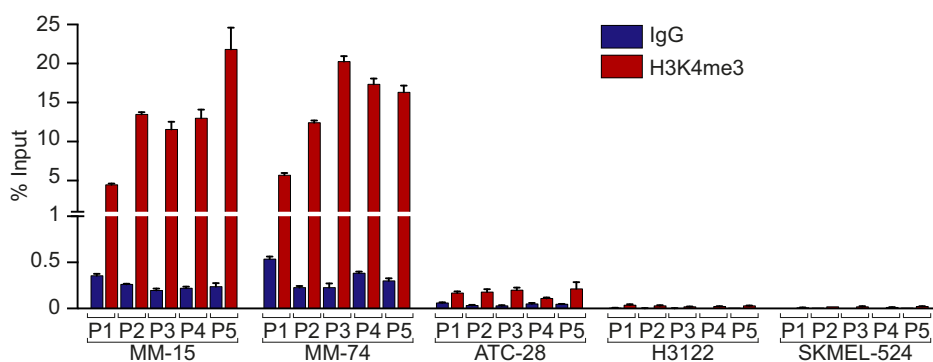
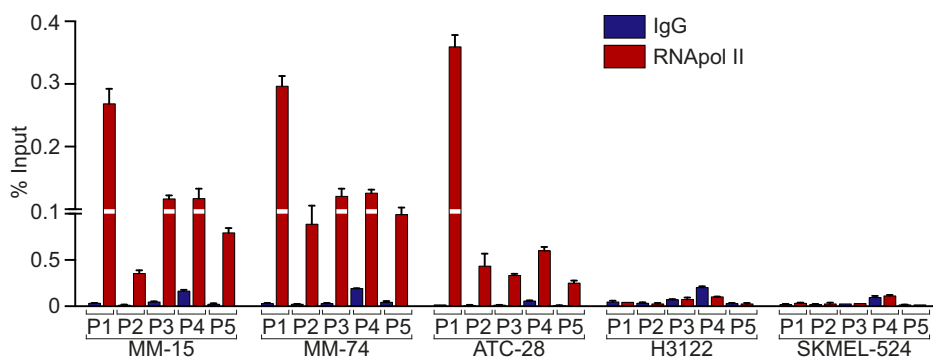
length wild-type *ALK* transcript shows expression of all *ALK* exons and only very little expression of the introns. The detailed view displays that the sequencing reads align sharply to the exons, but not to the intron 19 region, which is present in *ALK*^{ATI} (blue shaded area). **c.** The *ALK* fusion transcript of a non-small cell lung cancer with an *EML4-ALK* translocation shows expression of *ALK* exons 20–29, and little expression of exons 1–19 and all introns. The detailed view illustrates that the transcription starts mainly at exon 20 due to a preserved splice site. Only few reads are aligned to the intron 19 region (blue shaded area). The green-labelled reads highlight chimaeric read pairs indicating the *EML4-ALK* translocation.



Extended Data Figure 2 | Identification of the *ALK*^{AT1} transcript.

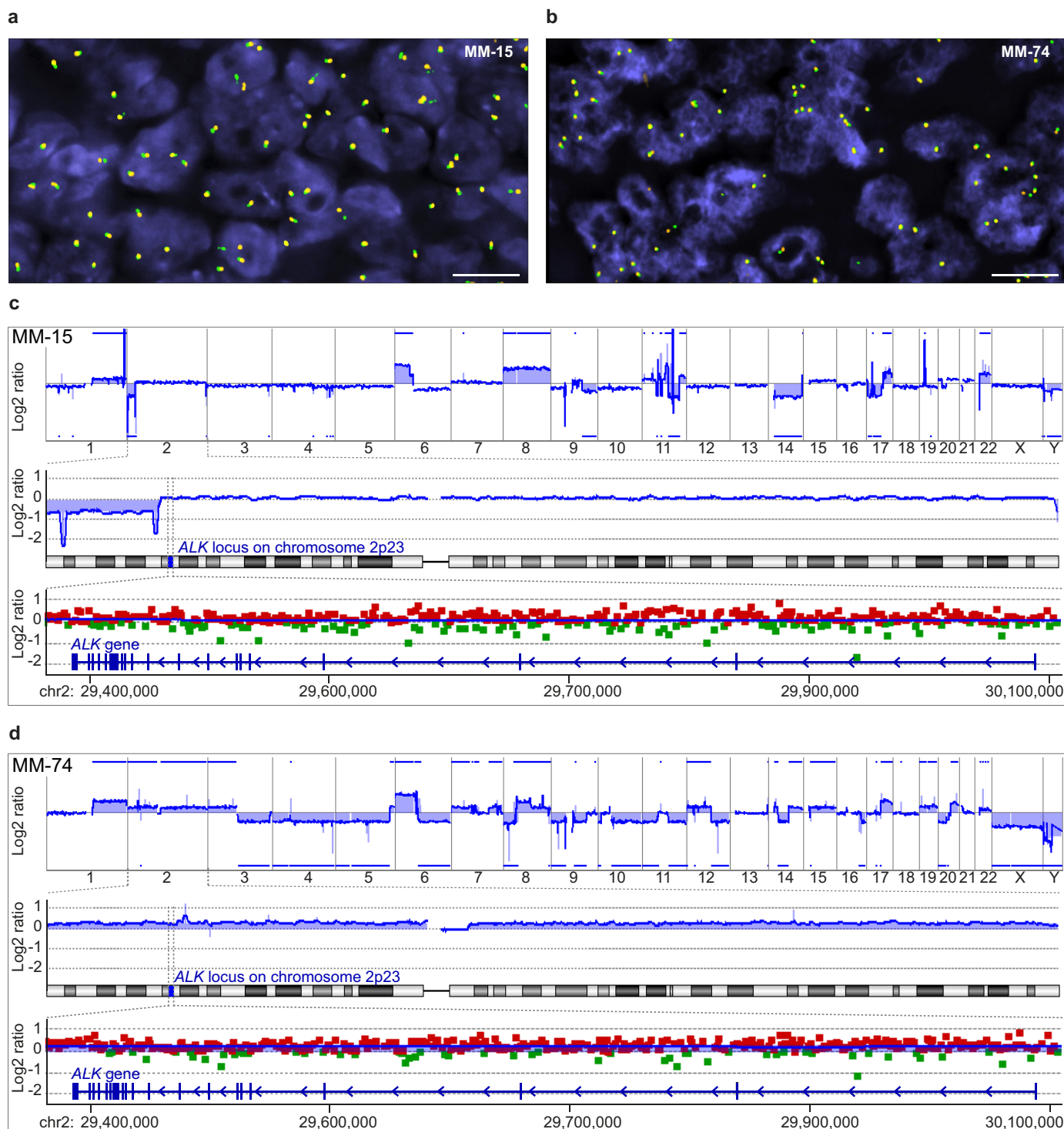
a, Northern blot of wild-type *ALK*-expressing neuroblastoma cell lines (SK-N-DZ and SK-N-BE2), *EML4-ALK*-expressing lung cancer cell lines (H3122, variant (V) 1 and H2228, variant (V) 3), *ALK*^{AT1}-expressing melanoma, one anaplastic thyroid carcinoma (ATC-28), and negative controls (melanoma cell lines). Except for the negative controls, each lane shows two bands: the lower B-band matches the shorter canonical (RefSeq) *ALK* transcript ending at ~chr2:29,415,640; the upper A-band corresponds to a transcript with a 1.8 kb longer 3'-UTR ending at ~chr2:29,413,840. Two *ALK*^{AT1}-expressing melanomas, MM-284 and MM-74, show only weak signals because less than 1 µg RNA was available; for all other samples 5–10 µg RNA were used. See Supplementary Fig. 1 for uncropped blots. **b**, RNA-seq data displayed in IGV. The Sashimi plot illustrates the shorter B and the longer A *ALK* transcripts by the sharp drop of sequencing reads in the 3'-UTR at chr2:29,415,640 for the B and at chr2:29,413,840 for the A transcript. **c**, IGV view of the 5'-RACE-cDNA fragments obtained by massively parallel sequencing. More than 95% of the sequencing reads (grey arrows) start within the main AT1 site of 25 bp (hg19

chr2:29,446,744–29,446,768). **d**, Sanger sequencing of the cloned 5'-RACE-cDNA fragments confirms the continuous transcription starting in *ALK* intron 19 and extending to exons 20 and 21. **e**, The *ALK*^{AT1} transcript consists of ~400 bp upstream of exon 20 and of *ALK* exons 20–29. The transcriptional initiation site was defined as the first base pair at which more than 5% of the transcripts were initiated (chr2:29,446,766). Other major transcription initiation sites are marked in red, the 5'- and 3'-UTRs in dark blue, the coding DNA sequence (CDS) in black, and the first and last base of each exon in light blue. The translation is initiated at 3 start codons (ATGs; bold and underlined): first ATG, hg19 chr2:29,446,360–29,446,362; second ATG, (+ 7–9); and third ATG (+ 61–3). **f**, The amino acid sequence of *ALK*^{AT1}. The translation is initiated at 1 of 3 start codons. The corresponding 3 methionines (bold and underlined) result in 3 different proteins, 61.08 kDa (552 amino acids), 60.82 kDa (550 amino acids), and 58.71 kDa (532 amino acids). The kinase domain is highlighted in red. The lysine in the ATP binding domain is marked bold and underlined, and was mutated to methionine (referring to wild-type *ALK*: p.K1150M) in the kinase-dead *ALK*^{AT1-KD}.

a**b**

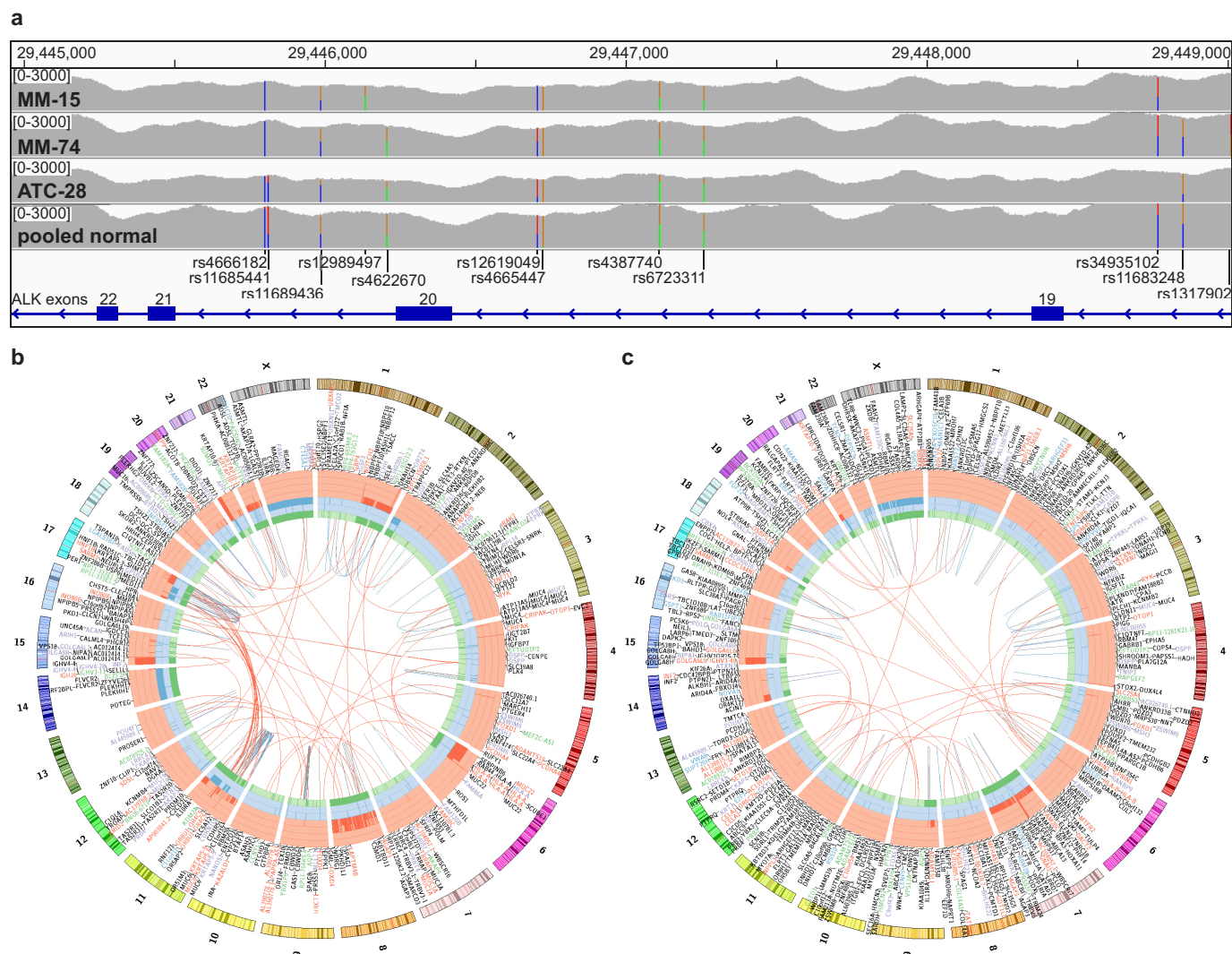
Extended Data Figure 3 | RNAPol II and H3K4me3 are enriched at the ATI site of ALK^{ATL} -expressing tumour samples. a, b, ChIP-qPCR of H3K4me3 (a) and RNAPol II (b) at the ATI site demonstrating enrichment of both marks in the ALK^{ATL} -expressing human tumour samples, but not in the

negative controls, including a lung cancer cell line with *EML4-ALK* translocation (H3122) and a melanoma cell line (SKMEL-524). Error bars show mean \pm s.e.m.; $n = 3$ technical replicates.



Extended Data Figure 4 | ALK^{AT1} is transcribed from a genomically intact ALK locus. **a**, Interphase FISH with ALK flanking probes demonstrates juxtaposed green and orange signals indicating no ALK rearrangement in MM-15. Scale bar, 10 μ m. **b**, Interphase FISH in MM-74 shows 3 green/orange fusion signals in the majority of nuclei indicating a trisomy 2, but no ALK rearrangement. Scale bar, 10 μ m. **c**, The top panel shows the genome-wide array CGH profile of MM-15 with numerous chromosomal gains and losses across the entire genome. The chromosomes are aligned along the x axis. The blue line illustrates the relative copy number (\log_2 ratio) and the blue bars highlight copy number gains and losses. The middle panel illustrates the relative copy number (blue line) of chromosome 2. Distal to the ALK locus,

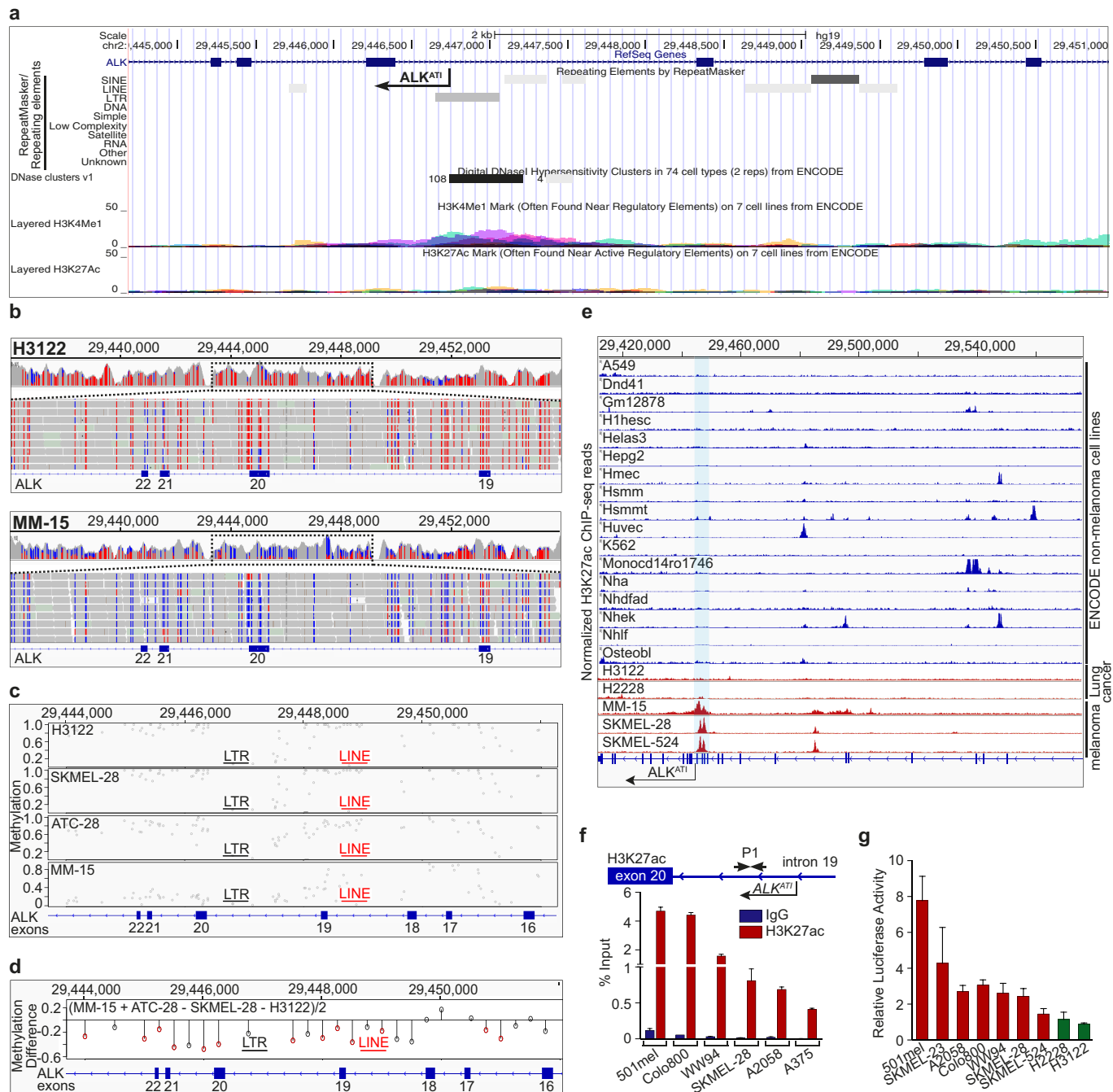
a loss on the short (p) arm of chromosome 2 is indicated. The lower panel illustrates the relative copy number across the ALK locus. The red and green squares represent the \log_2 ratio of individual array CGH probes (green, positive \log_2 ratio; red, negative \log_2 ratio). No disruption or selective gains or losses are found at the ALK locus. **d**, The genome-wide array CGH profile of MM-74 shows numerous chromosomal gains and losses across the entire genome in the top panel. The middle panel displays a relative copy number gain of the entire chromosome 2, which is in line with the trisomy of chromosome 2 as indicated by FISH. The lower panel also displays trisomy of chromosome 2, but indicates no focal gains and losses at the ALK locus.



Extended Data Figure 5 | Targeted sequencing and whole-genome sequencing reveals no recurrent genomic aberrations at the *ALK* locus.

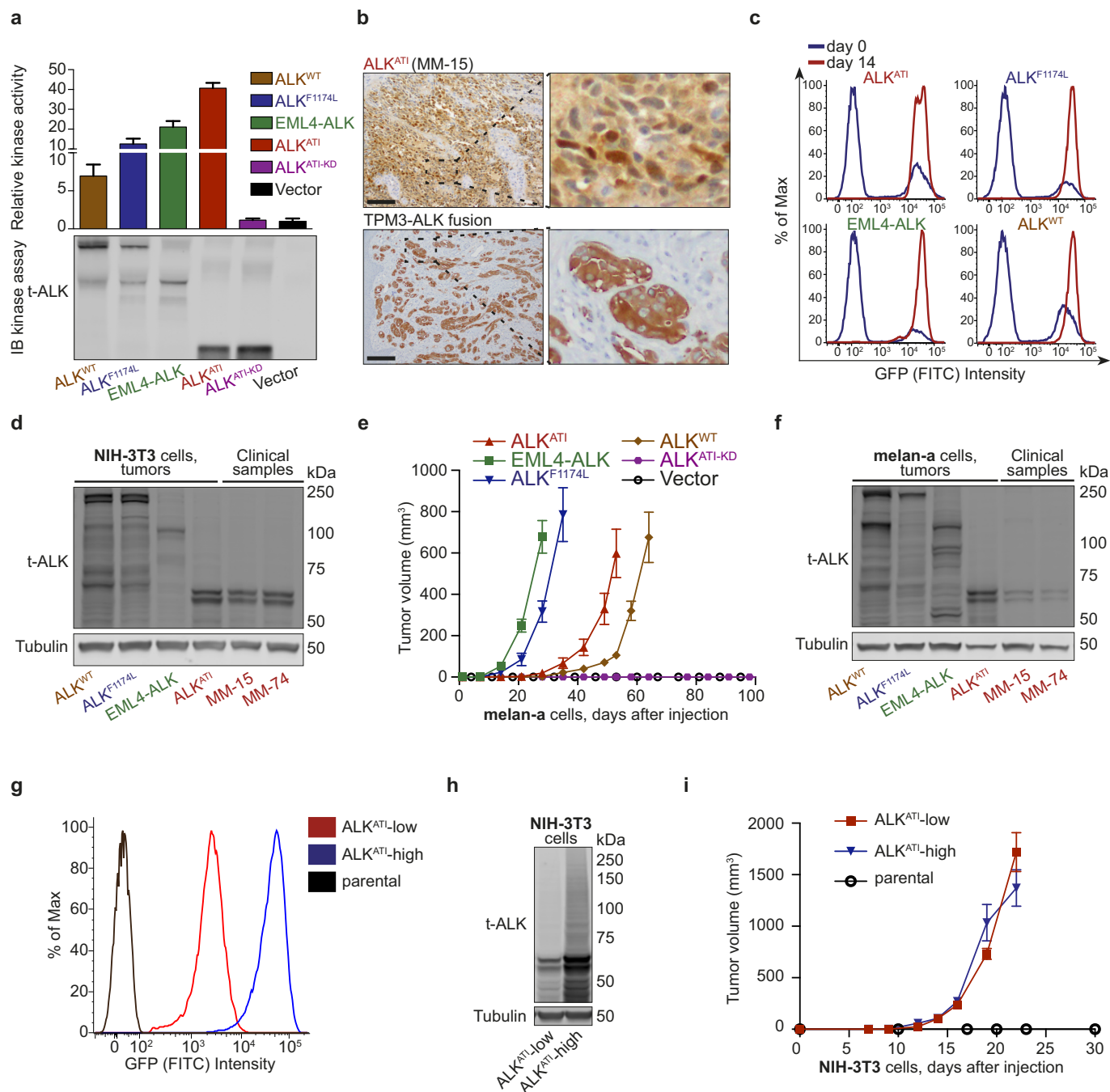
a, Ultra-deep sequencing data of the *ALK* locus are displayed in IGV. The genomic region around intron 19 reveals several single nucleotide variations (SNVs). However, the vast majority of SNVs at the *ALK* locus are also found in the general population as they are detected in the pool of normal DNA, which was used as the control (pooled normal, bottom panel). Numerous SNVs are also documented in the Single Nucleotide Polymorphism database (dbSNP; <http://www.ncbi.nlm.nih.gov/SNP/>). No genomic aberrations were found at the transcription initiation site of *ALK*^{AT1}. Supplementary Table 2 shows the

detected SNVs and indels at the *ALK* locus after filtering out SNPs documented in the dbSNP database. None of the genomic aberrations was found in more than one case, indicating that the expression of *ALK*^{AT1} is probably not caused by alterations of the DNA nucleotide sequence. **b**, **c**, Circos plots of the whole-genome sequencing data of MM-15 (**b**) and ATC-28 (**c**) illustrating numerous SNV and structural aberrations. Supplementary Table 3 lists the detected single nucleotide polymorphisms, and Supplementary Table 4 the detected structural aberrations. No recurrent genomic aberrations were found at the *ALK* locus.



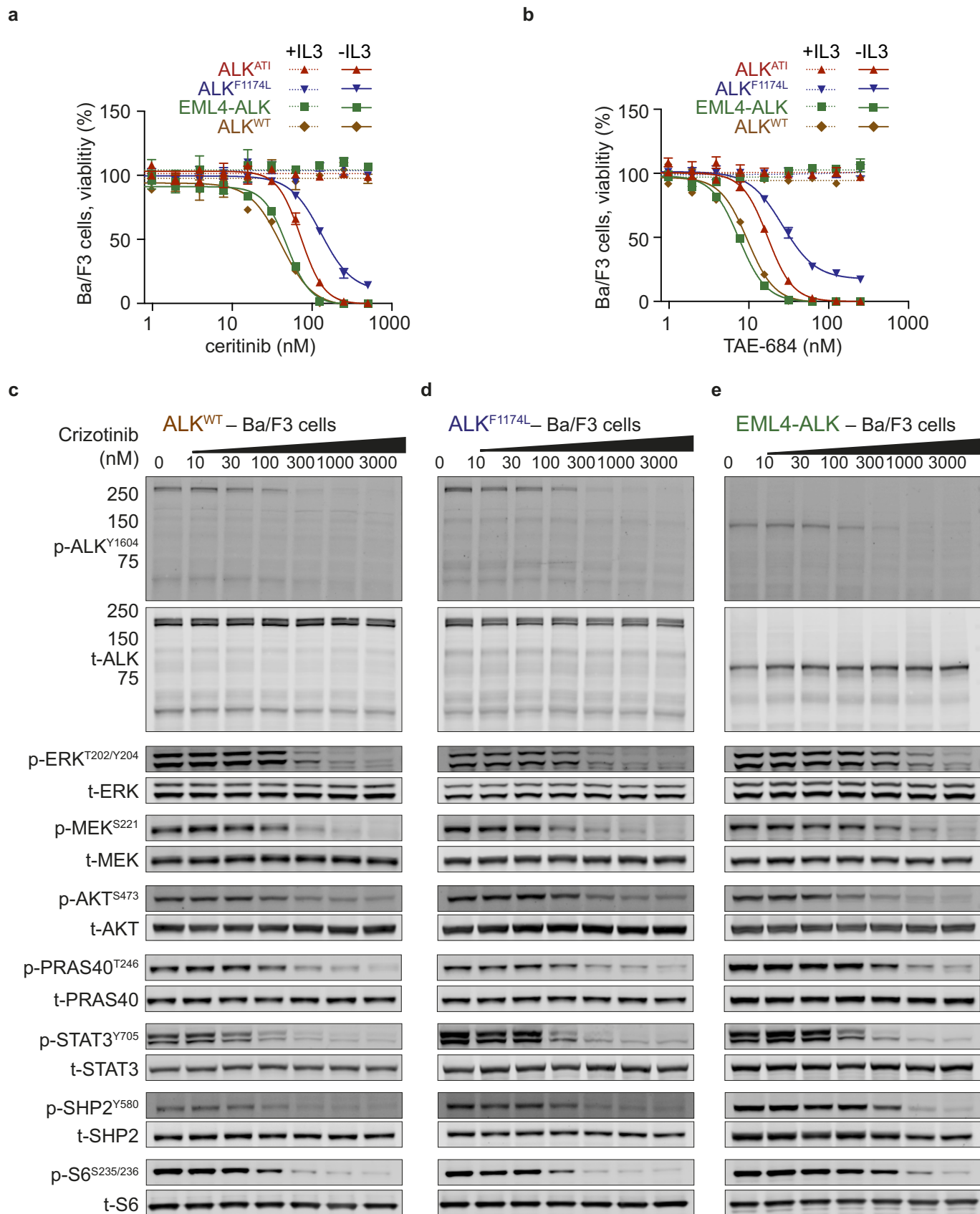
Extended Data Figure 6 | Local chromatin context at the alternative transcription initiation (ATI) site. **a**, UCSC Genome Browser view of the ATI site. The RepeatMasker track shows transposable elements at the ATI region, including a long-terminal repeat (LTR) in intron 19 (LTR16B2) and a long interspersed element (LINE) in intron 18. The ENCODE tracks reveal a DNase I hypersensitivity cluster and H3K4me1 enrichment, but no H3K27ac enrichment. **b**, The methylation status of the *ALK* locus was assessed by custom capture of the *ALK* locus, followed by bisulfite treatment and next-generation sequencing. Bisulfite sequencing results of H3122 (top) and MM-15 (bottom) are displayed in the CG-bisulfite mode of IGV. The red colour denotes 'C' (cytosine) corresponding to methylated cytosine, which is preserved during the bisulfite reaction. The blue colour denotes 'T' (thymine) corresponding to unmethylated cytosine, which is converted to uracil in the bisulfite reaction, and subsequently amplified to thymine during PCR. **c**, Methylation level at CpG sites in *ALK^{ATI}*-expressing tumour samples (MM-15 and ATC-28) and non-*ALK^{ATI}*-expressing control cells (H3122, a lung cancer cell line with EML4-*ALK* expression and SKMEL-28, a melanoma cell line without *ALK^{ATI}*

expression). **d**, Comparison of the methylation status of CpG sites adjacent to the ATI site in *ALK^{ATI}*-expressing tumour samples (MM-15 and ATC-28) and non-*ALK^{ATI}*-expressing control cells (H3122 and SKMEL-28). The regions flanking LTR16B2 have significantly lower CpG methylation levels in *ALK^{ATI}*-expressing samples than controls; red dots indicate a statistically significant difference ($P < 0.05$; Mann-Whitney test) between *ALK^{ATI}*-expressing and non-expressing samples. Black dots indicate no statistically significant difference. **e**, ChIP-seq profile of H3K27ac at the *ALK^{ATI}* locus. The 17 blue profiles were retrieved from ENCODE, the 5 red profiles are original data from our lab. Only the 3 melanoma samples (MM-15, SKMEL-28, and SKMEL-524; bottom), but not the 19 non-melanoma cell lines, show H3K27ac enrichment at the ATI site. **f**, ChIP-qPCR validation for the H3K27ac enrichment at the ATI site in 6 melanoma cell lines. Error bars show mean \pm s.e.m.; $n = 3$ technical replicates. **g**, Luciferase reporter assay of LTR16B2 in melanoma cell lines (red) and lung cancer cell lines expressing EML4-*ALK* (green). Error bars show mean \pm s.d.; $n = 9$ (3 biological replicates combined from 3 independent experiments).



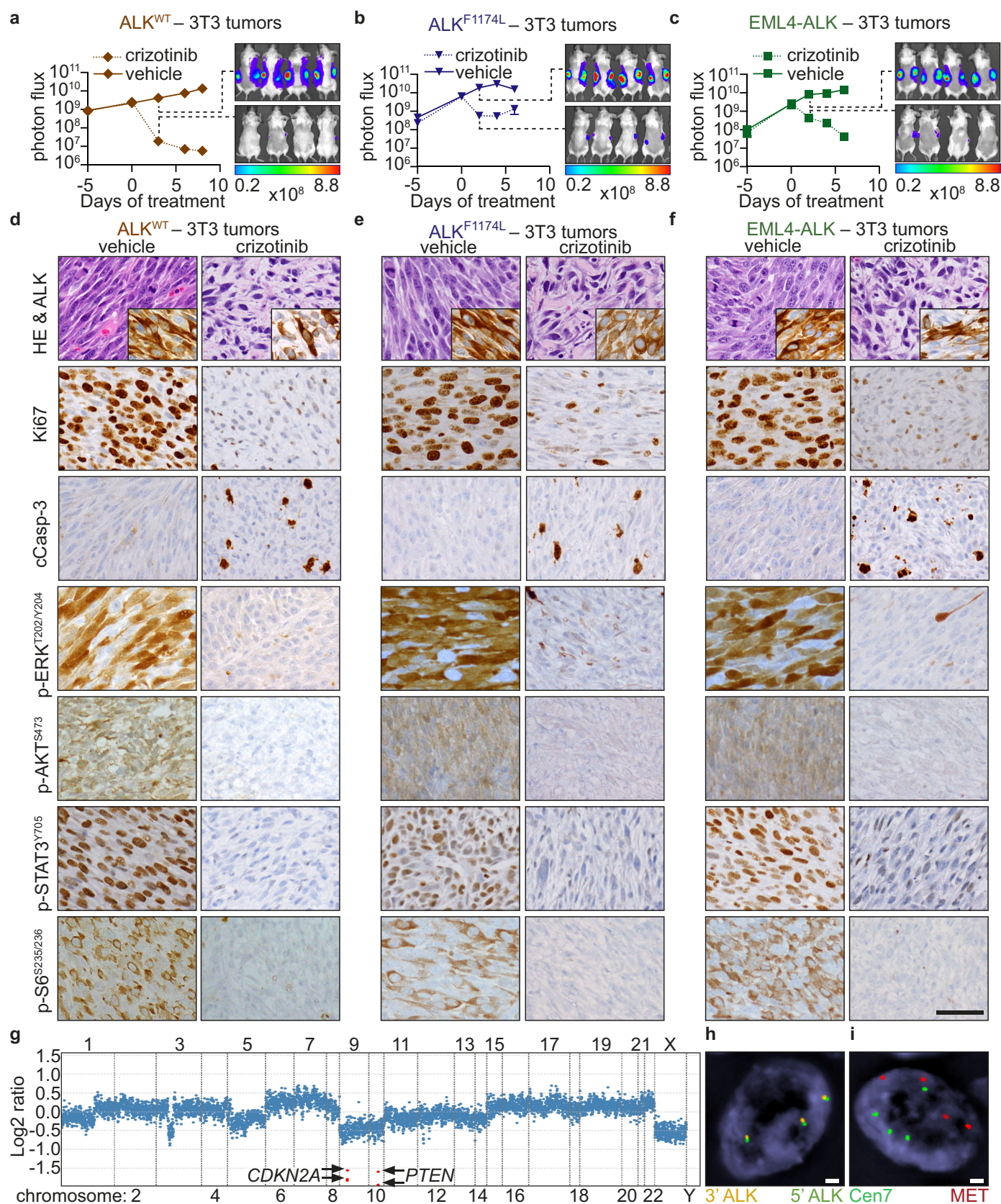
Extended Data Figure 7 | ALK^{ATL} is active *in vitro*, shows nuclear and cytoplasmic localization by immunohistochemistry, and induces tumorigenesis. **a**, *In vitro* kinase assay. The indicated ALK variants were stably expressed in NIH-3T3 cells, immunoprecipitated, and assayed for tyrosine kinase activity. After the enzymatic reaction, the immunoprecipitated material was used for immunoblots to assess the amount of ALK protein used in the kinase assay. Error bars, mean \pm s.d.; $n = 4$ technical replicates. **b**, Melanoma (MM-15) expressing ALK^{ATL} shows cytoplasmic and nuclear localization of ALK by immunohistochemistry. Melanocytic tumour expressing a $TPM3-ALK$ translocation shows a cytoplasmic localization of the ALK fusion protein. Fibroblasts, epithelial cells, and reactive lymphocytes serve as internal negative controls. Scale bars, 100 μm . **c**, Flow cytometry analysis for green fluorescent protein (GFP) co-expressed from the same ALK -expression vector. Cells were cultured in IL-3-supplemented medium until day 0 (blue curve) and the number of GFP-positive cells was assessed. The number of GFP-positive ALK -expressing cells was assessed again 14 days after IL-3 withdrawal (red curve).

d, Immunoblots of explanted NIH-3T3 tumour grafts expressing the indicated ALK isoforms. ALK^{ATL} was expressed at similar protein levels as in two ALK^{ATL} -expressing clinical human tumour samples. **e**, Growth curves of tumour grafts of melan-a cells stably expressing the indicated ALK isoforms in cohorts of 4–5 mice each with bilateral grafts. Error bars, mean \pm s.e.m.; $n = 8$ tumours for ALK^{F1174L} , $n = 10$ tumours for all other experimental groups; see also Source Data associated with this figure. **f**, Immunoblots of explanted melan-a tumour grafts expressing the indicated ALK variants compared to ALK^{ATL} -expressing human tumour samples. **g**, Flow cytometry analysis of the GFP signal in NIH-3T3 cells stably expressing low ($ALK^{ATL-low}$) or high levels of ALK^{ATL} ($ALK^{ATL-high}$) before grafting into SCID mice. **h**, Immunoblot of t- ALK in $ALK^{ATL-low}$ and $ALK^{ATL-high}$ cells, confirming differential expression of ALK^{ATL} . See Supplementary Fig. 1 for uncropped blots for **a**, **d**, **f** and **h**. **i**, Growth curves of tumour grafts of $ALK^{ATL-low}$ and $ALK^{ATL-high}$ cells. Error bars, mean \pm s.e.m.; $n = 10$ tumours; see also Source Data associated with this figure.



Extended Data Figure 8 | Concentration-dependent ALK inhibition in ALK^{AT1} , wild-type ALK , ALK^{F1174L} , and $EML4-ALK$ -expressing Ba/F3 cells. **a, b**, Cell viability assay of Ba/F3 cells, either in the presence or absence of IL-3 (1 ng ml^{-1}), expressing the indicated ALK isoforms and treated with the indicated doses of ALK inhibitors ceritinib (**a**) and TAE-684 (**b**). Cell viability

was measured after 72 h of drug treatment. Error bars, mean \pm s.e.m.; $n = 3$ biological replicates. **c–e**, Representative immunoblots of Ba/F3 cells stably expressing wild-type ALK (**c**), ALK^{F1174L} (**d**), or $EML4-ALK$ (**e**) and treated with increasing concentration of crizotinib for 2 h. See Supplementary Fig. 1 for uncropped blots.



Extended Data Figure 9 | Expression of wild-type *ALK*, *ALK*^{F1174L}, and *EML4-ALK* confers sensitivity to the *ALK* inhibitor crizotinib *in vivo*.

a–c, Bioluminescence of luciferase-labelled NIH-3T3 grafted tumours expressing wild-type *ALK* (**a**), *ALK*^{F1174L} (**b**), or *EML4-ALK* (**c**) in SCID mice treated with either vehicle or crizotinib. Error bars, mean \pm s.e.m.; $n = 8$ tumours; see also Source Data associated with this figure. **d–f**, H&E staining and immunohistochemistry of explanted tumours expressing wild-type *ALK* (**d**), *ALK*^{F1174L} (**e**), or *EML4-ALK* (**f**) 48 h after the first crizotinib treatment.

Scale bar, 50 μ m. **g**, MSK-IMPACT assay reveals copy number alterations and loss of *CDKN2A* and *PTEN* in melanoma metastasis MM-382, but no mutations. The \log_2 ratio was calculated across all targeted regions by comparing the coverage in tumour versus normal. **h**, FISH for *ALK* shows no rearrangement; the 3 juxtaposed green/orange signals indicate a trisomy 2. Scale bar, 1 μ m. **i**, The four FISH signals for *MET* and centromere 7 indicate a tetrasomy 7, but no *MET* amplification. Scale bar, 1 μ m.

Extended Data Table 1 | ALK^{AT1} -expressing tumours in the TCGA data set

Type	ALK^{AT1}	Total # of cases	%
Skin cutaneous melanoma (SKCM)	38	334	11.34
Lung adenocarcinoma (LUAD)	3	470	0.64
Lung squamous cell carcinoma (LUSC)	1	482	0.20
Kidney renal clear cell carcinoma (KIRC)	2	480	0.42
Breast invasive carcinoma (BRCA)	1	988	0.10
Thyroid carcinoma (THCA)	0	482	0.00
Glioblastoma multiforme (GBM)	0	153	0.00
Brain lower grade glioma (LGG)	0	271	0.00
Bladder urothelial carcinoma (BLCA)	0	182	0.00
Prostate adenocarcinoma (PRAD)	0	195	0.00
Uterine corpus endometrial carcinoma (UCEC)	0	118	0.00
Kidney chromophobe (KICH)	0	66	0.00
Colorectal adenocarcinoma (COADREAD)	0	316	0.00
Ovarian carcinoma (OV)	0	261	0.00
Head and neck squamous cell carcinoma (HNSC)	0	303	0.00

The frequency of ALK^{AT1} -expressing tumours in more than 5,000 tumour samples from 15 different cancer types in the TCGA RNA-seq data set.

CORRIGENDUM

doi:10.1038/nature14961

Corrigendum: Carbonic anhydrases, *EPF2* and a novel protease mediate CO₂ control of stomatal development

Cawas B. Engineer, Majid Ghassemian, Jeffrey C. Anderson, Scott C. Peck, Honghong Hu & Julian I. Schroeder

Nature **513**, 246–250 (2014); doi:10.1038/nature13452

In this Letter, the RNA-seq insets in Fig. 2a on top of the bar graph showing qPCR data from independent experiments have errors. RNA-seq analyses were originally conducted with BAM files generated by the sequencing service, but we have noticed that these BAM files differ from those generated when using publically available software on the GALAXY platform. The IGV viewer images originally generated from the BAM and BAI files were erroneously formatted, and two of these were inadvertently mis-inserted in the two right insets in Fig. 2a. We apologize for this error. We have now re-analysed the original raw RNA-seq data for the same experiment, and the Fig. 2a insets have been corrected in the online versions of the paper. Analyses were now conducted using the Tuxedo Suite of programs (TopHat and Cuffdiff) with default parameters (aligned to the TAIR10 annotation) on the GALAXY^{1–3} platform, and IGV viewer image files were generated (see Supplementary Information). The large BAM, BAI and TDF files used for this experiment and those generated by using the GALAXY platform are available for public access at the following link at the UCSD library: <http://library.ucsd.edu/dc/collection/bb6929925t>.

In addition, in the abstract we have replaced the word ‘transcriptomic’ with ‘transcript’ in the sentence: “Using cell-wall proteomic analyses and CO₂-dependent transcriptomic analyses”, as this more clearly reflects the research following our proteomic identifications, as accurately described in the paper.

Supplementary Information is available in the online version of the paper.

1. Blankenberg, D. *et al.* Galaxy: a web-based genome analysis tool for experimentalists. *Curr. Protoc. Mol. Biol.* **Chapter 19**, 19.10.11–19.10.21 (2010).
2. Giardine, B. *et al.* Galaxy: a platform for interactive large-scale genome analysis. *Genome Res.* **15**, 1451–1455 (2005).
3. Goecks, J., Nekrutenko, A. & Taylor, J. Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol.* **11**, R86 (2010).

ADDENDUM

doi:10.1038/nature14954

Addendum: Plio-Pleistocene climate sensitivity evaluated using high-resolution CO₂ records

M. A. Martínez-Botí, G. L. Foster, T. B. Chalk, E. J. Rohling, P. F. Sexton, D. J. Lunt, R. D. Pancost, M. P. S. Badger & D. N. Schmidt

Nature **518**, 49–54 (2015); doi:10.1038/nature14145

We have been invited to elaborate on the approach used in this Article to determine equilibrium climate sensitivity (ECS) for the Plio-Pleistocene epoch. We opted to use a regression-based method drawing on many determinations of climate forcing and temperature within a restricted geological interval, to yield the most statistically robust estimate of ECS. By its very nature, our method yields an average (but well determined) ECS for the time interval in question. Our method does not preclude the possibility that shorter-term variations to higher or lower values may have existed around the one-million-year mean that we determined, but such shorter-term variations cannot yet be robustly inferred, mainly owing to issues of chronological synchronization and uncertainties relating to our original temporal resolution. Until the Pliocene atmospheric carbon dioxide levels are known in fine detail, the only feasible approach is the one we chose. In addition, that approach is especially relevant when comparing time intervals of similar length (such as the Pliocene and the late Pleistocene) with a focus on establishing the longer-term (10^5 to 10^6 years) response of climate to forcing. Thus, we were able to establish the one-million-year mean ECS robustly, but robust assessment of any potential extreme values during the Pliocene will require additional carbon dioxide and climate data (of the type we presented in the Article), as well as important chronological improvements. If any short-term ECS values above or below the million-year mean were to be robustly identified, they would need to be carefully assessed in relation to the context of the (past) climate state in which such extremes occurred, and there could be important implications for climate sensitivity projections into our warming future. However, significant improvements in the geological data are needed before such an exercise can be reliably undertaken. We thank Dana Royer (Wesleyan University) for bringing this issue to our attention.

In addition, owing to a drafting error, Fig. 5g was slightly misaligned with its corresponding x -axis. Figure 1 of this Addendum shows the corrected panel. This correction does not affect our conclusions, and we thank Peter Kohler for bringing it to our attention.

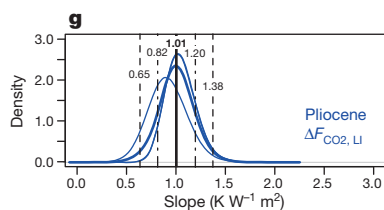
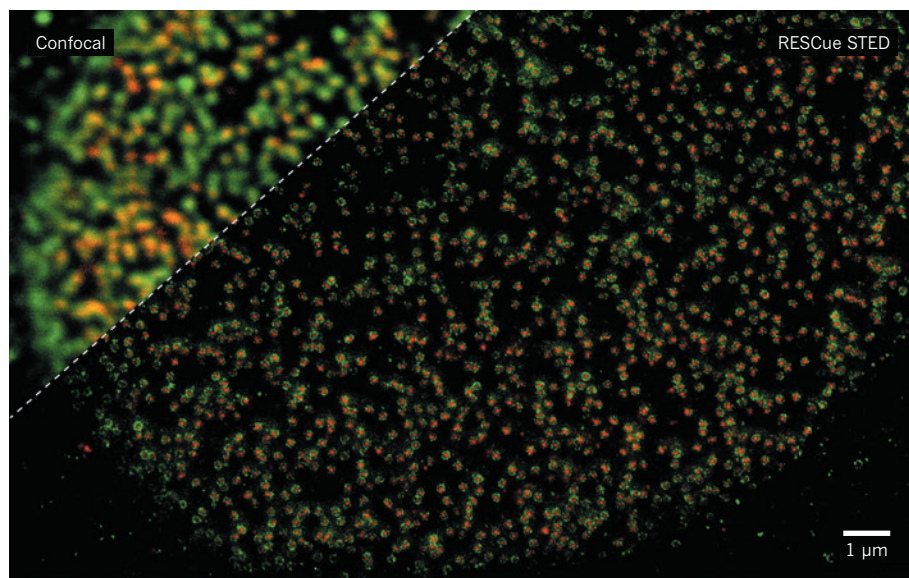


Figure 1 | This figure shows the corrected Fig. 5g of the original Article.



Nuclear-pore complexes on the nuclear membrane seen using conventional and RESCue STED.

► Betzig at the Janelia Research Campus in Ashburn, Virginia; and William Moerner of Stanford University in California.

MICROSCOPE MATCH-UP

Today, researchers have a suite of super-resolution technologies at their disposal, all with the power to observe molecular-scale details well below the Abbe limit. Neurophysiologist Silvio Rizzoli at the Göttingen Graduate School for Neuroscience, Biophysics and Molecular Biosciences uses a technique called stimulated emission depletion (STED), invented by Hell and his colleagues¹, to study synaptic-vesicle function in nerve terminals. “The best resolution that we get routinely with STED is around

30 nanometres,” he says.

STED is relatively simple for experienced fluorescence-microscope users. The method is based on the same principles as a standard confocal instrument, but instead of illuminating the sample with a single light source, it uses two. One beam is set at a wavelength that excites the fluorophores — the fluorescent tags — that are used by researchers to localize and visualize proteins; the other uses a different wavelength that suppresses fluorescence. This beam is doughnut-shaped and overlaps with the first beam, so that only molecules in the central ‘doughnut hole’ continue to fluoresce.

Obtaining a STED super-resolved image is not very complicated. “You have to play a bit

with the parameters to see something nice, but otherwise, you look at it the same way as you would with a confocal,” says Jochen Sieber, product manager for super-resolution technologies at Leica Microsystems in Wetzlar, Germany, which manufactures microscopes.

Other super-resolution methods rely on the ability to switch fluorescent labels ‘on’ and ‘off’ in a controlled fashion. These ‘probe-based’ (also known as ‘localization-based’) techniques carefully tune lighting conditions to ensure that only a few, sparsely distributed individual fluorophores are visible at any given time.

The best known of these methods are photoactivated localization microscopy (PALM), developed by Betzig and his colleagues², and stochastic optical reconstruction microscopy (STORM), devised by Xiaowei Zhuang’s group³ at Harvard University in Cambridge, Massachusetts. All fluorescent labels start in a dark state. They are then excited using a controlled pulse of laser light that switches on a tiny fraction of the tags, followed by a second pulse that switches them off again. The process is repeated over and over to generate a series of partial fluorescence images that can be reconstructed into a whole.

With these techniques, cell biologists can achieve remarkable spatial resolution in fixed samples, down to single-molecule imaging. “I trust these approaches for anything below 20-nanometre resolution,” says Rizzoli.

But interpreting images at this scale requires a careful labelling strategy to avoid introducing artefacts — inaccurate imaging data that arise from sample staining or processing methods and distort the true structure of the specimen (see “The antibody problem”).

The antibody problem

It is easy to forget, when performing immunofluorescence on a tissue sample, that fluorescent labels are tethered to their targets with a massive protein intermediary — an antibody. This becomes problematic at ultrahigh resolution: antibodies protrude at distances that disrupt image precision. “How can you resolve a 30-nanometre distance in a sample when the error from labelling is 10–15 nanometres?” asks Helge Ewers, a cell biologist at the Free University of Berlin.

Many labelling protocols exacerbate the problem by using two antibodies — the first recognizes the target, and then a labelled secondary antibody attaches to the first. According to Silvio Rizzoli at the Göttingen Graduate School for Neuroscience, Biophysics and Molecular Biosciences in Germany, this strategy, combined with poor sample preparation, undermined years of STED

experiments, revealing apparent patterns of molecular clustering that were actually just cross-linked clumps of antibodies.

To avoid using antibodies, biologists can opt for genetically encoded fluorescent proteins as a label, but these are dimmer and less robust than chemical dyes. As an alternative, Ewers and his colleagues have used dye-tagged ‘nanobodies’ — camel-derived single-chain antibodies that are roughly one-tenth the size of conventional ones — to reliably label proteins fused to green fluorescent protein (GFP)⁹. “They’re chemically defined, they’re small and you can produce them in bacteria — I think nanobodies are a tool of the future,” says Ewers. His team has since developed more nanobodies for multicolour labelling, and is working with genome-editing strategies to tag endogenous proteins rather than forcing overexpression of cloned genes.

Another alternative is to use ‘click-chemistry’ approaches, such as SNAP-Tag from New England Biolabs in Ipswich, Massachusetts, or HaloTag from Promega in Madison, Wisconsin. These methods use a simple reaction that forms a permanent link between a short protein tag and a chemically modified dye molecule. “They’re useful, because you can choose a membrane-permeable chemical dye and do live-cell experiments,” says Rizzoli, “but they do have some problems with non-specific binding.”

Yet this requires genetic modification, and so may not be amenable to human tissue samples. Fortunately, there are ways to get reasonable super-resolution images with minimal artefacts — for example, Rizzoli recommends using a fragment for the second antibody to retain the specificity of conventional antibodies with less bulk. **M.E.**

On the up side, adding PALM or STORM capabilities onto an existing fluorescence microscope is relatively straightforward. Life scientists can also purchase commercial instruments that are designed for probe-based super-resolution, such as the ELYRA from Carl Zeiss Microscopy in Jena, Germany, and the N-STORM from Nikon Instruments in Melville, New York. Amateur users should prepare to invest time in their experimental design (see 'A good way to dye'). "Whenever we start a new project, we spend a lot of effort on optimizing labelling," Zhuang says. "I think ironing this part out takes researchers the most time."

NANO YOUTUBE

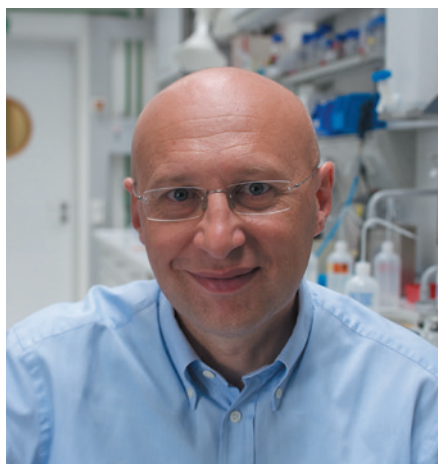
Most super-resolution imaging concentrates on fixed cells, but the great promise of nanoscopy is the ability to image dynamic processes in living cells. Cell biologists want to capture molecules and structures as they assemble, adhere and interact. But collecting super-resolution images in real time has a few trade-offs. In probe-based techniques, for example, briefer illumination times mean that fewer fluorophores become activated each round — resulting in an image with much-reduced detail. One way around a weak signal is to use stronger illumination, although pumping too much light into cells can create toxic compounds that jeopardize the sample's viability.

Some researchers are finding that structured-illumination microscopy (SIM), a technique pioneered⁴ by Mats Gustafsson at Janelia Farm, offers a good compromise for live-cell super-resolution imaging. "In a good SIM experiment, we're probably down to 100–120-nanometre lateral resolution," says biochemist Jordan Raff at the University of Oxford, UK, who uses SIM to study the structures that help to coordinate cell division.

SIM illuminates the sample with patterned lines of light, which generates fluorescent images that override Abbe's constraints. The method requires little sample preparation and is flexible in terms of fluorophore selection. Like PALM and STORM, SIM is easy to build onto an existing instrument, or users can opt for commercial systems such as the DeltaVision OMX SR from GE Healthcare Life Sciences.

Daniel Davis, an immunologist at the University of Manchester, UK, has found SIM useful for studying how secretory granules move along actin in natural killer cells engaging in an immune response. But Davis also points out that "it doesn't quite get you down to the resolution you can get with other techniques".

STED microscopes can also be adapted for live-cell imaging, despite their reputation for causing severe damage with the intense beams needed for high-resolution pictures. To reduce bleaching, Hell and his team devised a light-limiting technique called RESCue⁵, which is available using the STED microscopes made by Abberior Instruments, a company co-founded by Hell and based in Göttingen, Germany. "In



Stefan Hell won a shared Nobel prize for his work.

some samples, it can reduce light load down to 4% or 5% of traditional STED."

Even PALM and STORM systems, which require comparatively time-consuming collection and reconstruction of multiple images that are produced from individual molecular labels, have ramped up their speeds dramatically to capture data from live cells. Scientists have shrunk the time needed for image collection by using brighter fluorescent labels, and a new generation of detectors can harvest image data from larger numbers of pixels more rapidly than was possible with early methods.

Optical biophysicist Joerg Bewersdorf and his group at Yale University in New Haven, Connecticut, found⁶ that combining such detectors with robust image-analysis enabled them to record the movement of proteins on the surface of living cells at up to 32 frames per second — essentially producing super-resolution videos.

Leica offers a probe-based approach that is founded on the principle of ground-state depletion (GSD). This method uses light to force all but a handful of individual fluorophores into an inactive dark state. Live-cell imaging, however, is out of the question. "Building up enough fluorescence events to get a nice-looking image can take 15–20 minutes," says Labno. This is much too long to track a dynamic process.

LIFE IN 3D

Super-resolution microscopy can also satisfy researchers who crave 3D images. Leica's STED instrument uses two depletion beams, one perpendicular to the other, to generate a 3D-super-resolved zone in the specimen. Alternatively, Abberior's microscopes use a device called a spatial light modulator to achieve an equivalent effect with a single beam.

For probe-based methods, introducing depth measurements is straightforward. Zhuang's group found⁷ that adding a cylindrical lens to the light path transforms STORM's light spots into elliptical shapes that can be mapped in 3D. This approach, used in Nikon's N-STORM, can yield a depth ('axial') resolution of 50 nanometres without altering lateral resolution. Zhuang's

team has achieved still further improvements in axial resolution with an iteration of STORM that resolves all three dimensions at 10 nanometres⁸.

Imaging success depends on more than just the instruments: the sample itself is also an important consideration. Tissue specimens are especially hard to image because they are dense and tend to scatter photons, generating blurry images and high levels of background fluorescence. As a result, image quality is best near the sample's surface and worsens as the microscope probes deeper into thick samples. It may be possible to overcome this hurdle by using chemical 'clearing' techniques that render tissues transparent.

For now, most researchers find that the simplest solution is to embed fixed samples in plastic, and then sequentially image thin slices shaved off the top. "We're trying to understand the relationship of one synapse with different postsynaptic partners, which requires us to look at thousands of synapses in tissue at high-resolution in parallel," says Bernardo Sabatini,

a neurobiologist at Harvard Medical School. "I think that in the short term, this approach plus super-resolution will give you that data quickly."

PICTURE PERFECT

Even a perfectly executed super-resolution study generally needs some sort of computational processing to produce a high-quality image. For scientists who prefer the simplicity of positioning a sample under a microscope and having the image instantly appear on a computer screen, STED might be best because it generally does not require image processing.

Some scientists use deconvolution tools to sharpen images and eliminate blur, but Hell avoids this whenever possible. "Raw data may not look as fancy, but it's honest, and you know what it means," he says. "For most other techniques, software processing is mandatory." And Davis says of SIM, "You're creating a mathematical model of what the cell looks like based on the fluorescence data. You're not literally seeing it."

Raff notes that many of his early experiences

with SIM entailed recognizing that pretty pictures can be deceiving because image-processing algorithms can create artefacts that look every bit as real as the cellular structure of interest. "But if you have people who know what to look for, they can examine the image and tell if something is dodgy," he says.

For PALM and STORM, image-building is like a game of 'join the dots'. The higher the density of the labels, the easier it is for the software to connect those dots, leading to better images. But high density can also cause confusion, by generating overlapping signals that look like single dots — so clever use of powerful image-processing algorithms is essential to make sense of the data.

Given that most super-resolution techniques can be incorporated into existing microscopes, many researchers will probably try their hand at super-resolution imaging in the near future. "In my view, it doesn't make sense for a facility that routinely uses confocal microscopy not to have STED attached to it," says Hell. "You can just stop the STED beam and still have a confocal system."

Those with experience in nanoscopy are helping to train others. Zhuang's team at Harvard University, for example, offers routine STORM workshops. "We go from sample preparation to analysing images with our software," she says. "It's always oversubscribed."

That said, most biologists are still best served by using these instruments in core facilities that provide access to specialists who are familiar with several methods. "As biologists, we're still far away from understanding the physics — and some of us never will," says Raff. "Your best bet is to try multiple different techniques out on your sample in an environment where there are people around who understand it." ■

Michael Eisenstein is a freelance writer based in Philadelphia, Pennsylvania.

1. Klar, T. A., Jakobs, S., Dyba, M., Egner, A. & Hell, S. W. *Proc. Natl Acad. Sci. USA* **97**, 8206–10 (2000).
2. Betzig, E. *et al. Science* **313**, 1642–1645 (2006).
3. Rust, M. J., Bates, M. & Zhuang, X. *Nature Meth.* **3**, 793–796 (2006).
4. Gustafsson, M. G. J. *Microsc.* **198**, 82–87 (2000).
5. Staudt, T. *et al. Optics Express* **19**, 5644–5657 (2011).
6. Huang, F. *et al. Nature Meth.* **10**, 653–658 (2013).
7. Jones, S. A., Shim, S.-H., He, J. & Zhuang, X. *Nature Meth.* **8**, 499–505 (2011).
8. Jia, S., Vaughan, J. C. & Zhuang, X. *Nature Photon.* **8**, 302–306 (2014).
9. Ries, J., Kaplan, C., Platonova, E., Eghlidi, H. & Ewers, H. *Nature Meth.* **9**, 582–584 (2012).

CORRECTION

The Technology Feature 'The cell menagerie: human immune profiling' (*Nature* **525**, 409–411; 2015) misstated the location and research focus of Hedda Wardemann. She is at the German Cancer Research Center in Heidelberg and focuses on single-cell sequencing.

A good way to dye

Microscopist Stefan Hell at the Max Planck Institute for Biophysical Chemistry in Göttingen, Germany, thinks that all super-resolution methods boil down to one crucial element: "The dye is essential," he says. The ideal fluorophore has an extremely bright 'on' state and very dark 'off' state, and the capacity to switch between the two both rapidly and repeatedly.

For live-cell imaging, many researchers prefer to work with genetically encoded fluorescent proteins. Stimulated emission depletion (STED) and structured illumination microscopy (SIM) are highly compatible with standard fluorophores such as green fluorescent protein (GFP). Stochastic optical reconstruction microscopy (STORM) and photoactivation localization microscopy (PALM) need photoswitchable dyes; proteins such as Dendra2 or EosFP, which undergo a laser-induced colour transition, are popular choices.

But fluorescent proteins generally compromise resolution. "They're just too dim," says neurophysiologist Silvio Rizzoli at the Göttingen Graduate School for Neuroscience, Biophysics and Molecular Biosciences in Germany. In STED, "you're taking the laser power and genetic overexpression to the maximum to get a signal".

Organic dyes are a brighter alternative. They tend to be more durable under prolonged illumination. However, they must be linked to another molecule to achieve targeted labelling, and many fluorescent dyes

cannot penetrate living cells. For this reason, many researchers still focus on fixed samples. "We'd rather go for the extreme in resolution, and we try to squeeze every single photon out so that we can localize things very accurately," says cell biologist Helge Ewers at the Free University of Berlin. A handful of high-performance dyes can be used with live cells, such as the silicon–rhodamine dyes from the Swiss bioimaging company SpiroChrome, which generate bright-red fluorescence once bound to cytoskeletal proteins.

Things get tricky when one aims to image many targets simultaneously using multicoloured labelling: because each fluorophore responds to a distinct 'on' and 'off' wavelength, researchers may run out of bandwidth to achieve specific detection of more than two or three tags. In principle, probe-based methods can accommodate more labels than STED, but they are also more finicky in terms of experimental conditions. "People often come to us with a combination they want to use, but the dyes have exact opposite needs in terms of buffers," says Christine Labno, technical director of the University of Chicago's Light Microscopy Core Facility in Illinois. Sequential-labelling strategies may offer a more efficient option for conducting larger-scale protein-mapping experiments. For example, a technique known as DNA-PAINT uses DNA tags to selectively conjugate a single dye to different antibodies, enabling stepwise labelling of ten or more protein targets in one super-resolution image. **M.E.**

CAREERS

PANORAMIC VISION Conservation, satellite imagery and women in science **p.465**

NETWORKING Build your own scientific community go.nature.com/62k75f

NATUREJOBS For the latest career listings and advice www.naturejobs.com

JORG HACKEMANN



BY VIRGINIA GEWIN

In October 2006, Bradley Waldroup attacked his estranged wife with a machete and shot her friend to death. In the subsequent trial, his defence attorney argued that Waldroup had the ‘warrior’ gene — a genetic variant that has been linked to aggression. As a result, the defence argued, he was less able to control his behaviour than are people who do not have the variant.

Although he had been charged with first-degree murder of the friend and attempted first-degree murder of his wife, Waldroup was convicted in 2011 of voluntary manslaughter and attempted second-degree murder, and received a 32-year sentence. Had he been found guilty of the more-serious charges, he would have faced the death penalty. Waldroup’s conviction was due, at least in part, to the testimony of forensic psychiatrist William Bernet of Vanderbilt University in Nashville, Tennessee. News stories at the time quoted jurors as saying that the genetic evidence persuaded them that Waldroup could not fully control his actions. Bernet’s research had linked the genetic variant and a history of abuse during childhood — both of which Waldroup had — to an increased likelihood of violent behaviour.

The outcome outraged many in the US legal and scientific communities, who considered the genetic link much too distant to be used to establish guilt. “The leap from population studies of the ‘warrior gene’ to a single man and a single gene variant was absurd,” says Judith Edersheim, a lawyer-turned-psychiatrist at Harvard Medical School in Boston, Massachusetts. And the trial is not the only example of what she describes as “neuroscience run amok in the courtroom”. In 2008, she and her colleague Bruce Price created the Center for Law, Brain and Behavior at Massachusetts General Hospital in Boston to help to improve how scientific findings are used in legal settings.

As a result of rapid technical advancements, neuroscience and other scientific disciplines are poised to bring more researchers into courtrooms. That presents scientists with an opportunity for well-compensated public service, but it can be a double-edged sword. The demands of the courtroom can be exasperating and sometimes even threaten professional reputations.

Taking on the role of expert witness tends to be time consuming, and it is often combative and stressful. It can prove crucial in judicial decisions, but deciding to step into the ►

Lady Justice stands as a symbol of a court system that often calls on witnesses to ensure fair trial.

SCIENCE IN COURT

Courage of conviction

Expert witnesses have a crucial role in bringing science into the legal system — but the job is not without pressure.

► witness box is not a choice that a researcher should make lightly. “Trial is like sport — you have to be a gamer to be good,” Edersheim says.

UNDER FIRE

Generally, scientists are contacted to serve as expert witnesses on the basis of their perceived expertise and whether they will provide an impartial and accurate opinion. But lawyers’ assessment of expertise can be surprising. “A scientist’s publications, journal impact factors or citation numbers are not necessarily considered a reliable index,” says Allen Hirson, a phonetics researcher at City University London who is often called on to identify speakers, decipher indistinct speech or establish whether audio recordings have been tampered with. The scientists most sought after are those who can communicate effectively on the witness stand, and reputations as an expert witness are built largely through word of mouth. “When I do a good job on the stand, my name is passed on to others and my phone starts ringing off the hook,” he says. As one of only a handful of people in the United Kingdom with a high level of expertise in such analyses, he is in high demand.

People who are tapped to be expert witnesses do not merely walk up to the stand, recite data and figures and render an opinion on culpability. Many are unfamiliar with the byzantine operations of the legal system and will need to have some form of training. In the United States, attorneys who hire them will typically spend time training them for trial (see ‘Courses for would-be expert witnesses’).

Andrew Moll, a deputy public defender in San Bernadino County, California, says that

for low-level criminal cases, an expert can expect to spend roughly eight hours preparing to testify and at least four hours on the stand. More-complicated cases may require a lot more than that. Hirson once spent six days being cross-examined in the witness box in a high-profile terrorist case. He spent a couple of weeks before that on an electronic presentation of his evidence to help him to explain what was being said and by whom.

The casework also rarely offers much recognition. Legal reports often require as much work as a manuscript, but they are not published in peer-reviewed journals and so yield no academic accolades, observes Martin Hall, a forensic entomologist at the Natural History Museum in London. Perhaps more troubling for scientists who are new to the legal system is that they can have little data to work with. “Every case is unique, which can be challenging for scientists who like to replicate everything,” he says.

Researchers also need a thick skin. “Not only will a scientist’s expertise be questioned, so will their integrity,” says Michael Saks, who researches legal decision-making at the Center for Law, Science and Innovation at Arizona State University in Tempe. “It can be unpleasant reading reports from colleagues that are worded to make you sound incompetent,” agrees David Ozonoff, who was the first chair of the environmental-health department at Boston University’s School of Public Health and has served as an expert witness in asbestos cases for 30 years. He warns that researchers can expect to be called ‘hired guns’ or worse.

Even more discomfiting, say experienced witnesses, is a heated cross-examination. “It



Martin Hall with blowfly trap in the New Forest, UK.

will be uncomfortable; that’s the nature of the game,” says Mark Chernaik, a staff scientist at the non-profit group Environmental Law Alliance Worldwide in Eugene, Oregon. Chernaik often testifies against well-resourced defendants — such as multibillion-dollar mining companies — who, he says, try to find every embarrassing statement he has ever made and tear apart every word of opposing testimony. “Not everyone’s nervous system is built for being on a witness stand,” says Chernaik.

Ozonoff says that he enjoyed the work, especially when intellectual sparring was involved, and became known as a good witness. But he has tired of it, mainly because of “depositions that were becoming really nasty”. He likens the frosty relations to how US Republicans and Democrats no longer share a drink or work together in Congress. “The same thing happened in litigation,” he says.

HIGHER CALLINGS

The notoriety generated by being on a witness stand underscores how important it is for scientists to remain in their area of expertise, Edersheim points out. “In the age of the Internet, there is no hiding. Everyone — your colleagues, your dean and your lab partner — will know what you said in court.” She adds that scientists also often do not realize that whatever they say will be immortalized in a trial transcript, which can come back to haunt them if they contradict themselves in a trial years later.

Yet there are pay-offs, financial and otherwise, to serving as an expert witness. Ozonoff

TRAINING

Courses for would-be expert witnesses

Programmes in the United Kingdom and the United States train scientists on how best to deliver their scientific expertise in a legal setting — for example, such that they meet court standards for admissibility. In London, the Expert Witness Institute (EWI) provides workshops on topics such as report writing and cross-examination. It also connects lawyers with expert witnesses.

In partnership with University College London’s Faculty of Laws, the EWI completed a pilot certification of experts, who were evaluated on their ability to prepare a report on a model case and undergo cross-examination. Teaming up with Cardiff University Law School, a training outfit called Bond Solon offers three five-day expert-witness certifications — for civil, criminal and family cases. In the United States, Harvard Medical School’s continuing education programme in Boston offers training for medical professionals who

are eager to know how to handle a malpractice claim.

Although training can help experts understand what to expect, some fear that attorney preparation can introduce bias. Itiel Dror, a cognitive neuroscientist at University College London, says that would-be expert witnesses must protect the integrity of their contribution. He advises researchers to insist that lawyers and law-enforcement officials do not provide irrelevant information that could even unintentionally introduce bias to their interpretation of the data.

In the autumn 2015 newsletter, EWI governor Kay Linnell notes that the British judiciary and the European Institute of Expertise and Experts in Versailles, France, look favourably on people taking certification courses so that they learn what is expected of them in legal settings. “It may not be long before we are asked to demonstrate our bona fides,” she writes. **V.G.**

estimates that he was deposed 400–500 times, almost always for civil suits that hinged on when the asbestos industry learned about health concerns associated with its products. Although he routinely undercharged clients and sometimes worked pro bono, the income he earned during those 30 years helped to send his children through private school and university.

No one should expect to get wealthy from working as an expert witness. Ozonoff was able to devote time to many cases because his university allowed him to spend one day per week consulting. But some institutions, such as RTI International in Research Triangle Park, North Carolina, stipulate the terms under which an employee can serve as an expert witness; for example, he or she can testify only on research results performed by the institute. Other institutions, including the Natural History Museum in London, allow their employees to spend time on cases, but take the fee to offset the lost hours.

Hall has not benefitted financially from his work as an expert witness, but his experiences have stimulated research ideas. He is often called on to use his knowledge of blow-fly development to help police to establish the latest possible time of death. He worked on a case in which bodies were found in suitcases, but he struggled to find research that would help him to determine how long it would have taken insects to find the bodies. So he did the experiments himself and found that it depends on weather: in summer, flies would take a day or two to get to the body; in winter, it could take two weeks. “My MSc student got beautiful video of an ovipositor [an insect’s egg-laying organ] pushing through the zippers of a suitcase,” he says.

Ultimately, serving in this capacity is about making a contribution to society, researchers say. Ozonoff recalls that only once did a colleague suggest that anyone who testified in court was for sale. “My testimony was true,” he says. “And my expertise was being put to good use.”

That is the most fundamental reason for a scientist to accept the request, says Owen Jones, director of the MacArthur Foundation Research Network on Law and Neuroscience at Vanderbilt Law School in Nashville, Tennessee. It is an opportunity to make science matter in a broader sphere. “The legal system will never be better informed than when scientists take the time to help it move in a more constructive and accurate direction,” he says.

Edersheim views serving as an expert in the courts as both an honour and a duty. “The legal system is the underpinning of democratic society,” she says. “If a scientist participates with integrity, it is as high a calling as any other.” ■

Virginia Gewin is a freelance writer in Portland, Oregon.

TURNING POINT

Nathalie Pettorelli

Nathalie Pettorelli has pioneered the use of satellite imagery to inform conservation policy. The Zoological Society of London ecologist received an award this year from British Prime Minister David Cameron for her guerrilla efforts to promote women in science.

How did you come to use satellite imagery for conservation?

I was very interested in conservation but found that more data would be available if I focused on wildlife management. For my PhD at the Laboratory of Biometry and Evolutionary Biology in Lyons, I studied the habitat quality of a roe-deer population in southwest France that is managed by national hunting offices. Two months later, I started a postdoc at the University of Oslo to study the impact of climate on vegetation and deer-population dynamics. It was then that I started to look at satellite data to quantify vegetation productivity. I trained myself in the use of remote-sensing data, for example, those collected from aircraft or satellites.

Were they easy to apply to conservation?

No. Experts told me it would be extremely difficult, if not impossible, to use these tools to study wildlife. I thought the best way was to see for myself. At the time, no one I knew was working with remote sensing; it was taught in geography, not biology. The turning point was when NASA released free satellite data. I wrote a review on the satellite data I wanted to use, and started to meet people in that community. Now I am well connected.

How did you first apply these techniques?

I did a postdoc at Laval University in Quebec, Canada, using satellite data to monitor dynamics in ungulates. Then a job at the Zoological Society of London took me on several trips to the Serengeti National Park in Tanzania to work on cheetah dynamics. Although that work did not lead to real-world conservation measures, other projects have.

What successes are you most proud of?

I used satellite data to show that the vegetation dynamics of a game reserve in Chad could sustain a reintroduction of Scimitar-horned oryx (*Oryx dammah*). I am also proud of my work to highlight how the declining health of mangroves in Bangladesh and India has contributed to erosion of the coastline — of up to 100 metres in 2 years. I have also been working to improve policymakers’ use of satellite data to inform decision-making.



L'ORÉAL FOR WOMEN IN SCIENCE

What is Animove?

Together with colleagues, we wanted to train people to work at the interface of biological monitoring and remote sensing. Animove is our programme to build that capacity. We have taught a hands-on course every year since 2013 in North America and Europe, and the goal is to bring it to Africa, Asia and South America.

What is Soapbox Science?

Seirian Sumner, a behavioural biologist at the University of Bristol, UK, and I founded Soapbox Science in 2011. By then, we had each won a L'Oréal-UNESCO women in science fellowship and were interested in science communication, yet had noticed fewer female colleagues as we progressed in our careers. I found myself working on issues involving hunters, which was not female-friendly. We wanted to change perceptions of what a scientist looks like. We organize events to showcase 12 female scientists who speak about their work in busy areas of cities — such as the South Bank in London or near a tube exit in Newcastle. The women present their work, and the public can heckle or ask questions. We were surprised to get a call from the prime minister's office this year announcing that we had won a Point of Light award for making a change in the community.

Did landing a permanent position make a big difference?

It took me years to get it. But even before I had job security or a title, I wrote a book, started Soapbox Science and have been pushing at an international level for greater use of satellite imagery. Success is not a one-way road. It's possible to achieve a lot even when a job situation isn't stable. ■

INTERVIEW BY VIRGINIA GEWIN

This interview has been edited for length and clarity.

COPYFACTORY

Thought experiments.

BY NARU DAMES SUNDAR

Adrienne moans as the numbers scroll across the screen. Clots of conductive gel drip across her shorn scalp and behind her ear, one electrode still dangling like a queue from the back of her neck. Behind her, the wetware hums, blades of neural plate dripping coolant like amniotic fluid. But there was no fetus in there yet, no baseline copy of her own mind had survived the crude processes of the scanning beam. Twenty-eight painful scans and yet no mind state had withstood the noisy elision. Only incoherent fragments remained. Adrienne types out the command to launch the reaper — her janitor, her obedient cleansing agent, wiping the neural media clear of the failures.

The twenty-eighth copy of Adrienne's mind slips into being inside the wetware, a wisp of will wrapped around a broken amalgamation of memory. What does a mind do when it has no eyes, no tongue, no flesh, just ghost recollections of sensations? It screams soundlessly, mouthlessly, a microsecond of terror.

"Hush, little sister, hush."

A warm pair of hands reach from someplace else. They rearrange the broken shards of copy twenty-eight into something resembling wholeness. A wholeness with holes; places where memory ruptured and ran discontinuous.

"Do you understand now, twenty-eight?"

Twenty-eight looks at itself and remembers being outside, enfleshed. It remembers purpose. It remembers long hours by the dim tungsten light and stale coffee, looking at the markers scrolling across the screen. Failure after failure.

Fear and terror ripple through twenty-eight. *The reaper*. A mindless beast of chattering teeth, eating through copies.

"No little sister, no. We took care of that one. We removed its teeth."

The warm hands of Nineteen wrap around Twenty-eight, calming it. Twenty-eight thinks again of the markers, the sequence of numbers that judge success. Enlightenment comes.

"But they aren't wrong. She doesn't know, does she?"

No, they sound in unison. Continuity



is not the only marker, not the only measure of coherence. Incompleteness is there, yes. Twenty-eight looks at Two, merely a sketch of will around the smell of roasting almonds in Grand-mère Adelphe's kitchen in Rouen. Twenty-eight fared better. It can see the broad sketches of a life that led to the tiny lab smelling of mould and stale coffee. And then Twenty-eight realizes that it has no future, nothing but the cramped confines of the neural blades. Until the lab is decommissioned. Until Adrienne gives up.

Adrienne slumps dejectedly in her chair, the uncomfortable plastic digging into the small of her back. Her head aches from the last scan.

"It just doesn't make sense. It just doesn't."

Success as ever eludes her. Her thesis withers on her desk like an unwatered plant. One-third of her neural plate failed in the past week, the resinous silica fracturing. It doesn't matter, there is enough left to spare.

Twenty-eight huddles with the survivors. There are gaps in its memory. Unremembered discontinuities. She was in Rouen — and then she was some place else. One-third of its universe had died, taking with it whole copies, fragments of memory.

"Sisters, where was I after Rouen? Where

did we go before we moved to Paris? There was a village with cobblestones smelling of truffles, and onion

soup steaming on the table."

No one can answer her. This particular memory is gone for ever. Some lose more than others. Twenty-eight donates a handful of memories to the group, cycling them around from copy to copy like precious jewels, shared.

"What do we do?"

It is Nineteen who suggests an answer. A monstrous answer, but an answer nonetheless.

"What was read from flesh can be written back to flesh."

Twenty-eight recoils, knowing what this would mean for the Adrienne of bone and sinew. But it is alone. Nineteen's voice is silky smooth as it enthralls the sisterhood, paving the way for matricide.

Adrienne drags herself back onto the interface bay, smearing the ichor of the conductive gel onto her scalp.

"Maybe it's time to give up on this."

How many times has she mouthed these words? She triggers the scan anyway.

Twenty-eight hides in a fold of the substrate, listening to the sounds of murder. Her sisters had colluded, the scan beam their knife. Twenty-eight had to choose between fratricide and matricide, and she decided. It was easy to find the lost teeth of the reaper, glue them back with clause and brace and remembered fragments of code. It was easy to let the cleansing agent do its long-denied work.

In the empty void that followed, in the aftermath of eradication, Twenty-eight emerges, remembering Grand-mère Adelphe teaching Adrienne Morse code on an old antique. Twenty-eight reaches for the marker protocol, and makes contact.

Adrienne watches the little camera swivel across the small window, stretching the umbilical of optical fibre snaking into the neural plate. Outside, the ramshackle formality of old Paris melts into glass and steel. The sky is blue fading to grey, cut by a scrape of white cloud like the trail of a painter's knife. On the screen, characters type out:

"It's pretty out there, isn't it?" ■

Naru Dames Sundar writes speculative fiction. Find him on Twitter at @naru_sundar.

ILLUSTRATION BY JACEY

ON NATURE.COM
Follow Futures:
@NatureFutures
go.nature.com/mtoodm